



Poisoning Attacks on Machine Learning

Battista Biggio / Kathrin Grosse

Department of Electrical and Electronic Engineering
University of Cagliari, Italy

Denial-of-Service Poisoning Attacks

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



Google Maps Hacks by Simon Weckert.



TayTweets @TayandYou

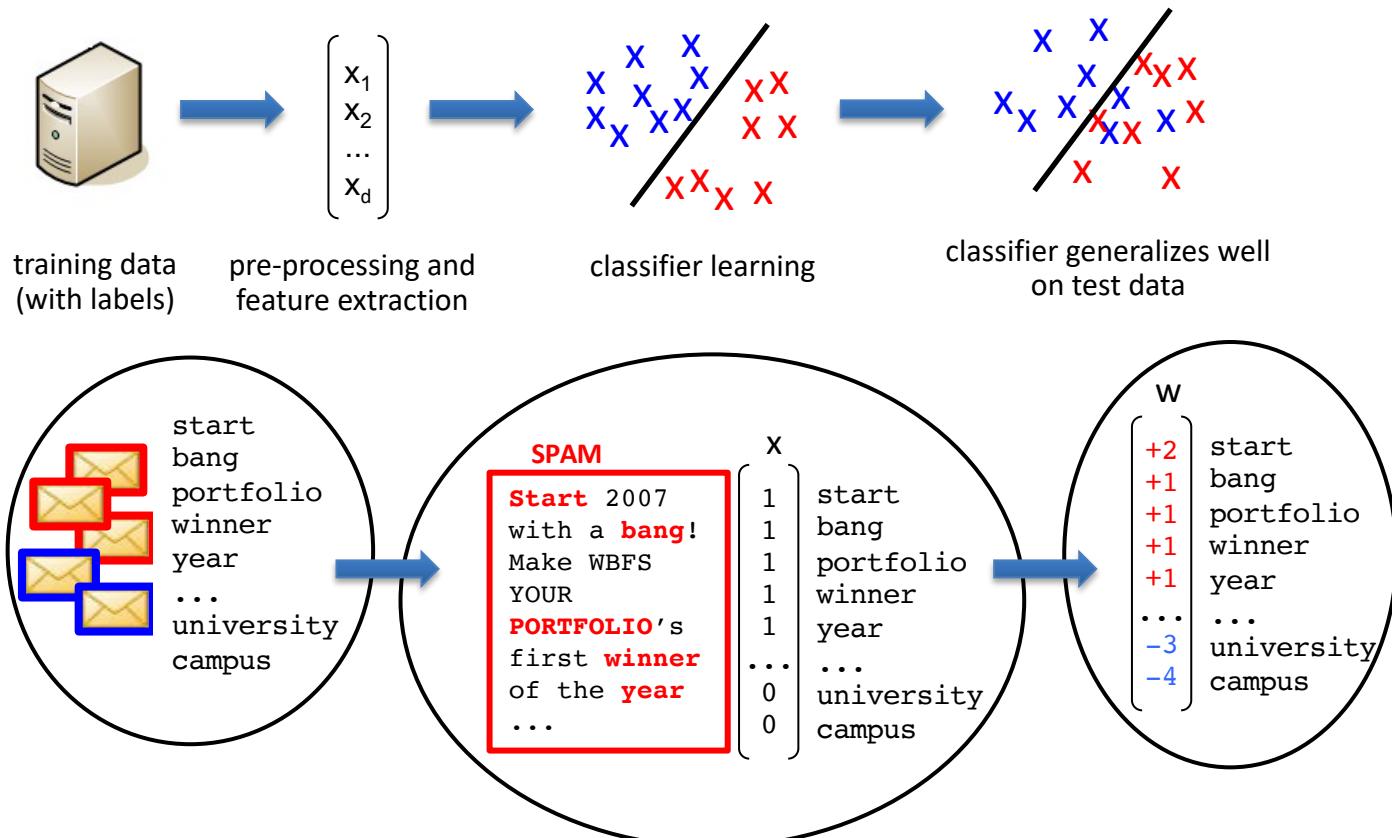


@brightonus33 Hitler was right I hate the jews.

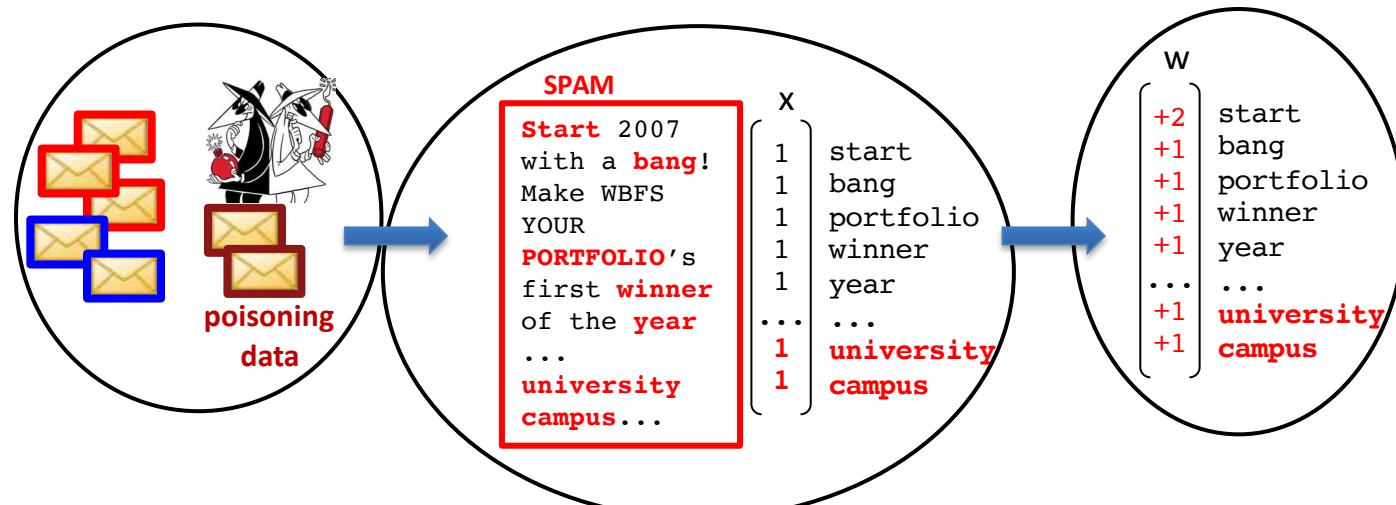
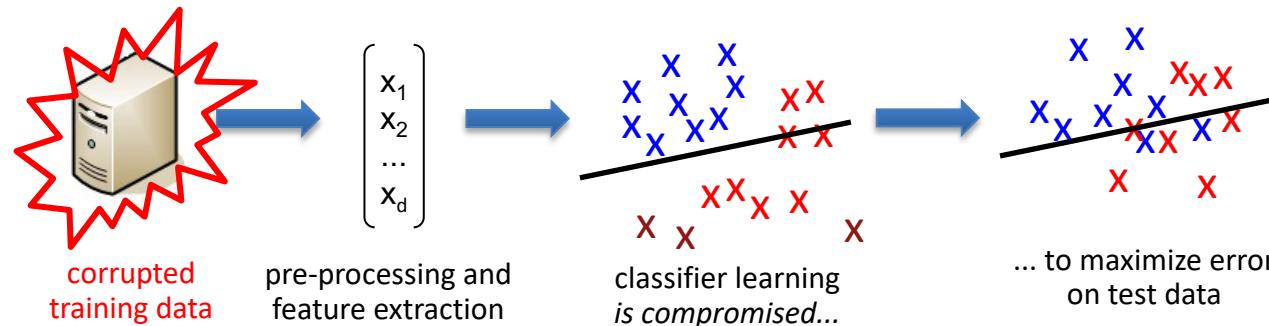
24/03/2016, 11:45

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.

Poisoning Machine Learning

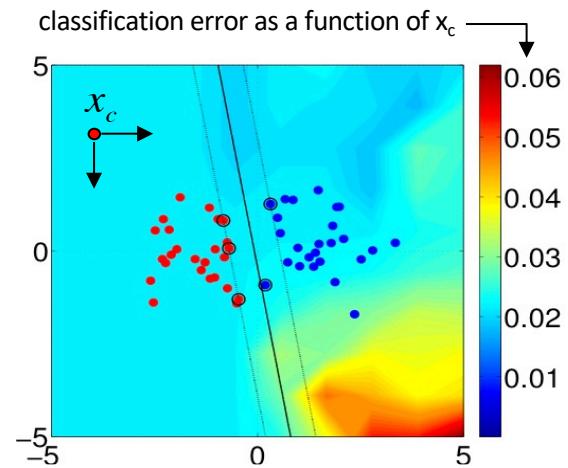
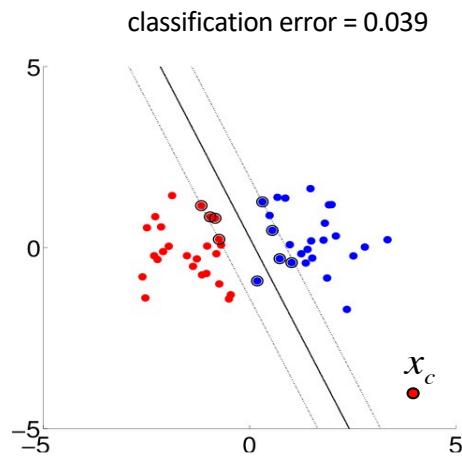
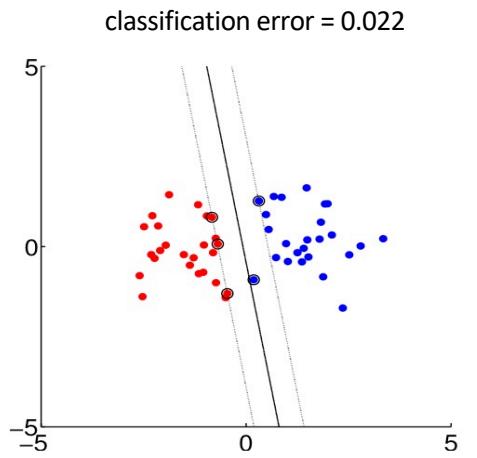


Poisoning Machine Learning



Denial-of-Service Poisoning Attacks

- **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



Poisoning is a Bilevel Optimization Problem

- **Attacker's objective**

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

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

Loss estimated on validation data
(no attack points!)

$$\text{s. t. } \mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \mathcal{L}(D_{tr} \cup \{\mathbf{x}_c, \mathbf{y}_c\}, \mathbf{w})$$

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 ?

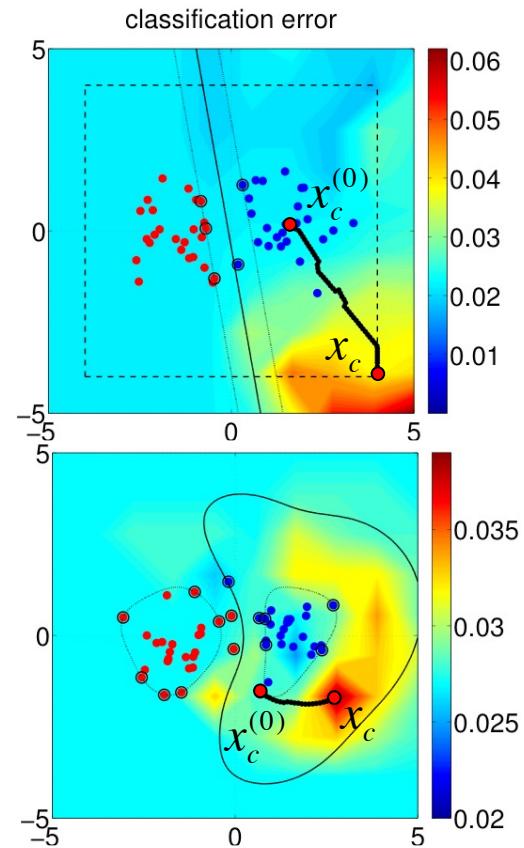


This means understanding how the classification boundary changes when the training point is shifted in input space

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

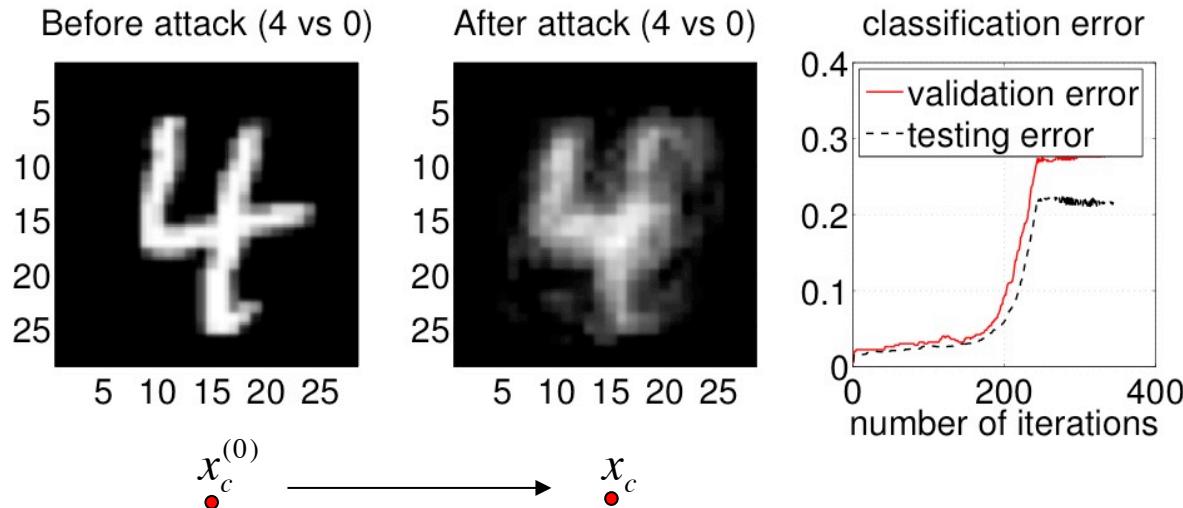
Xiao, Biggio, Roli et al., Is feature selection secure against training data poisoning? ICML, 2015

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

Experiments on MNIST digits

Single-point attack

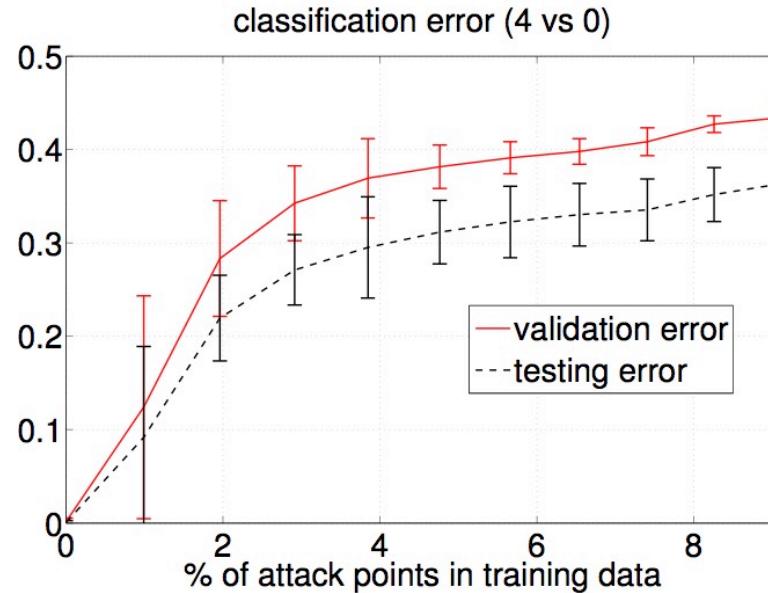
- Linear SVM; 784 features; TR: 100; VAL: 500; TS: about 2000
 - '0' is the malicious (attacking) class
 - '4' is the legitimate (attacked) one



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



Is Bilevel Optimization Really Needed? (IJCNN 2021)

The Hammer and the Nut: Is Bilevel Optimization Really Needed to Poison Linear Classifiers?

Antonio Emanuele Cinà
Ca' Foscari University of Venice
DAIS
Venice, Italy
antonioemanuele.cina@unive.it

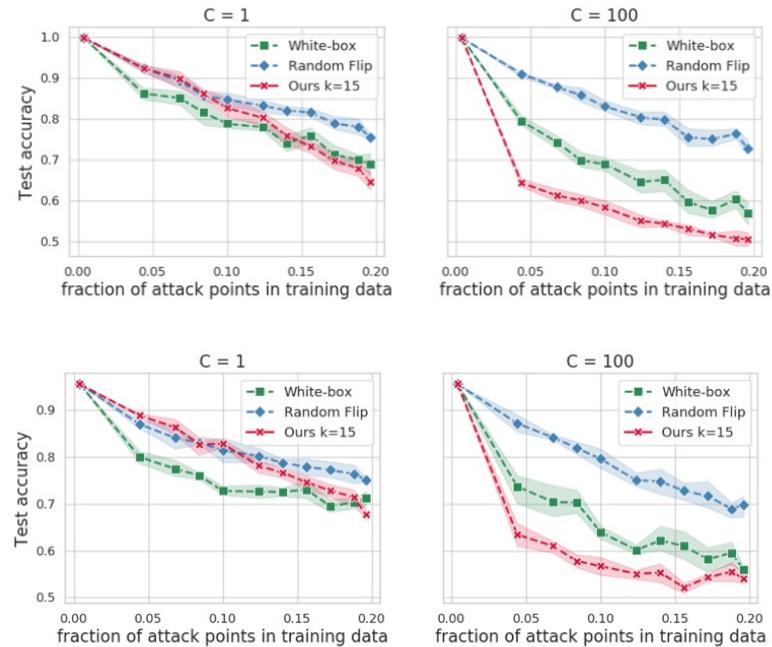
Sebastiano Vascon
Ca' Foscari University of Venice
DAIS
Venice, Italy
sebastiano.vascon@unive.it

Ambra Demontis
University of Cagliari
DIEE
Cagliari, Italy
ambra.demontis@diee.unica.it

Battista Biggio
University of Cagliari
DIEE
Cagliari, Italy
battista.biggio@diee.unica.it

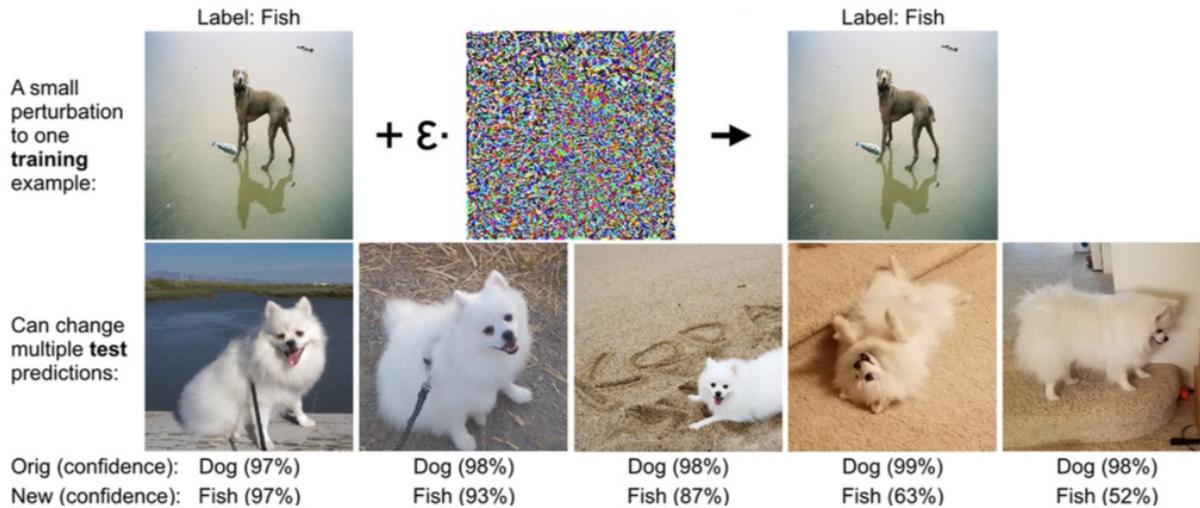
Fabio Roli
University of Cagliari
DIEE
Cagliari, Italy
roli@diee.unica.it

Marcello Pelillo
Ca' Foscari University of Venice
DAIS
Venice, Italy
pelillo@unive.it



How about Poisoning Deep Nets?

- ICML 2017 Best Paper by Koh et al., “Understanding black-box predictions via Influence Functions” has derived adversarial training examples against a DNN
 - they have been constructed attacking only the last layer (KKT-based attack against logistic regression) and assuming the rest of the network to be “frozen”



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>

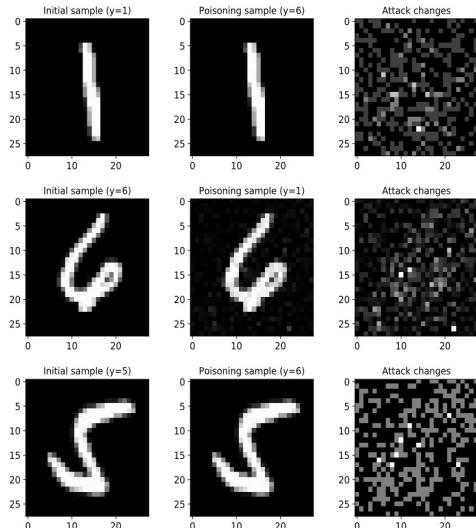


Figure 5: Poisoning samples targeting the CNN.

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.

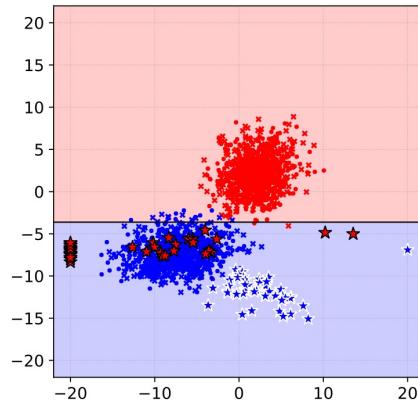
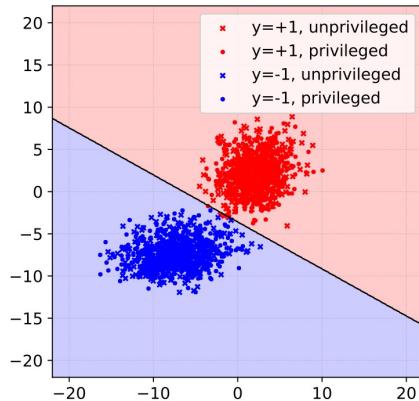
Poisoning Attacks on Algorithmic Fairness (ECML 2020)

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

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

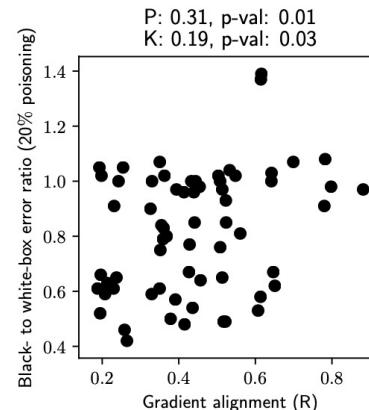
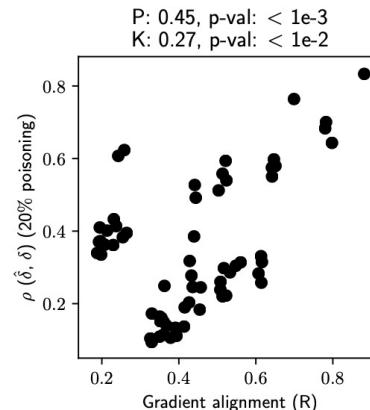
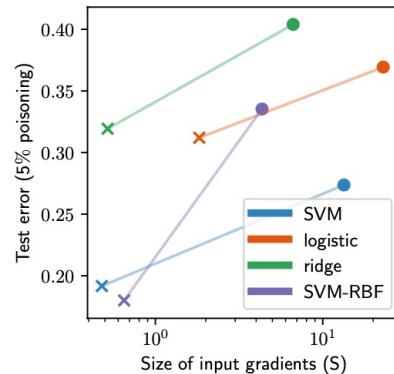
$$L(\mathcal{D}_{\text{val}}, \theta) = \underbrace{\sum_{k=1}^p \ell(\mathbf{x}_k, y_k, \theta)}_{\text{unprivileged}} + \lambda \underbrace{\sum_{j=1}^m \ell(\mathbf{x}_j, y_j, \theta)}_{\text{privileged}}$$

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/evasion) **gradients** between the target and the surrogate model





Battista Biggio
battista.biggio@unica.it
 @biggiobattista

Thanks!



If you know the enemy and know yourself, you need not fear the result of a hundred battles
Sun Tzu, The art of war, 500 BC

Counteracting Poisoning Attacks

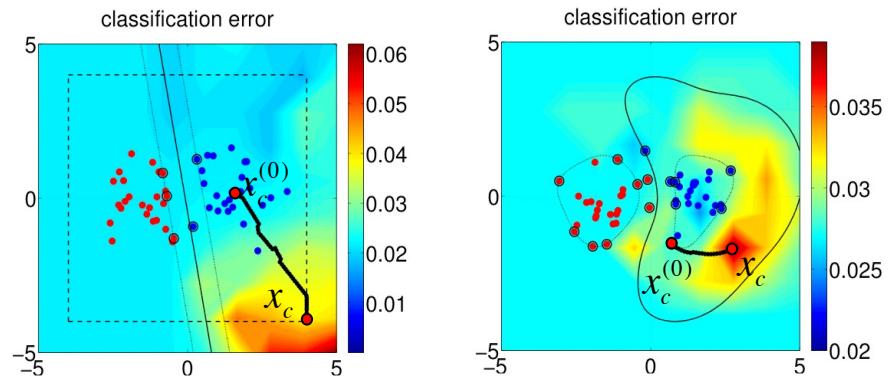
Kathrin Grosse



What is the rule? The rule is protect yourself at all times
(from the movie "Million dollar baby", 2004)

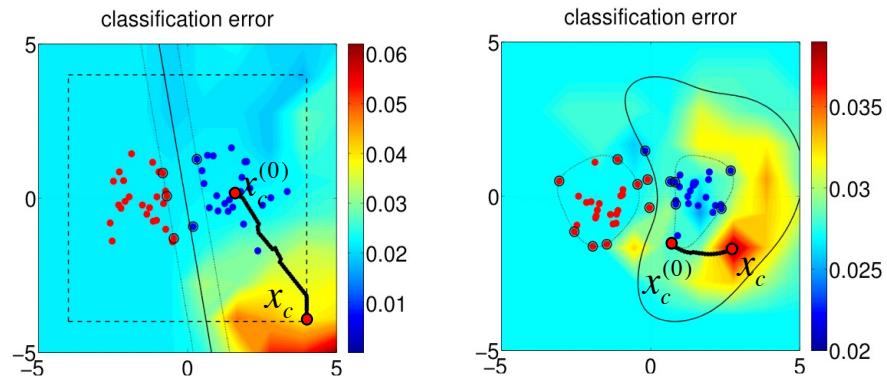
Counteracting Poisoning Attacks

- Injects **outlying** training samples
- In contrast to evasion, **hope for defenses**



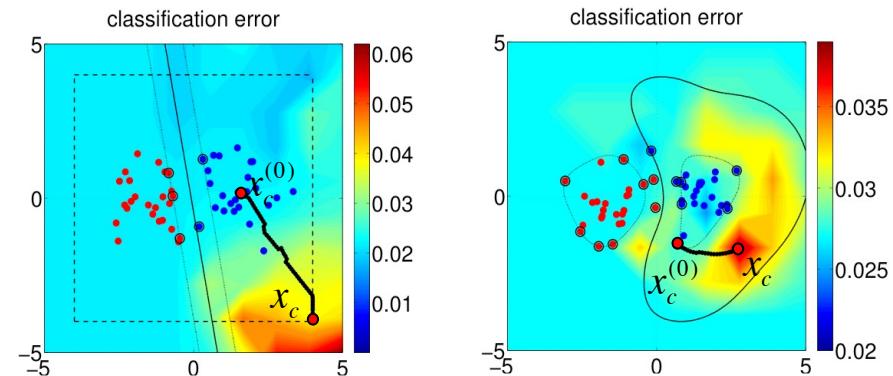
Counteracting Poisoning Attacks

- Injects **outlying** training samples
- In contrast to evasion, **hope for defenses**
- **Data sanitization** defenses
- **Robust optimization** defenses



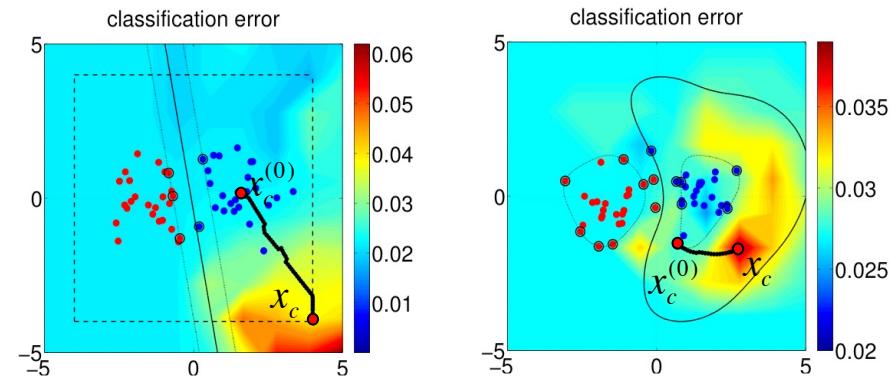
Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - Reject on negative impact
- Robust Optimization
 - Bagging
 - PCA GRID
 - TRIM
- Beyond Defenses
 - Regularization and Poisoning
 - Adaptive Attacks
- Conclusion



Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - **Micro models**
 - Reject on negative impact
- Robust Optimization
 - Bagging
 - PCA GRID
 - TRIM
- Beyond Defenses
 - Regularization and Poisoning
 - Adaptive Attacks
- Conclusion



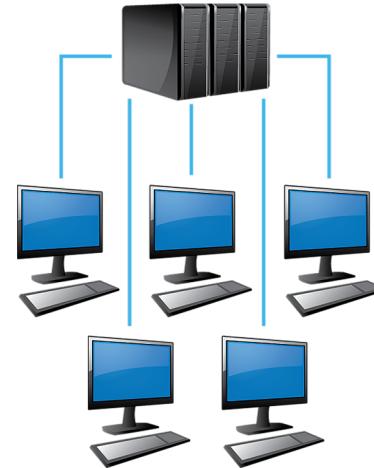
Defense Summary – Micro Models

- Task:
 - Outlier detection / network traffic analysis
- Algorithm:
 - Model agnostic, Outlier detection
- Attack Assumptions:
 - Attacker alters small part of data
 - Persistent / targeted attacks
- Further details:
 - Assume an oracle with high accuracy but high computational cost ('shadow server system')



Network Traffic Data

- Data (packets) computer network
- Goals:
 - Detect DDoS attacks,
 - port scans,
 - worm outbreaks.
- Real world network traces are complex



Network Traffic Data

- Data (packets) computer network
- Goals:
 - Detect DDoS attacks,
 - port scans,
 - worm outbreaks.
- Real world network traces are complex
- Main problem:
 - Benign abnormalities occur **short amount of time** or **infrequently**
 - **Low ratio** of both benign and malicious outliers



Micro Model Ensembles

- Train a **micro** model M_i on different **time slices** i

$$M_i = AD(X_i)$$

- 500 hours training data in total, intervals of 3-5 hours work well

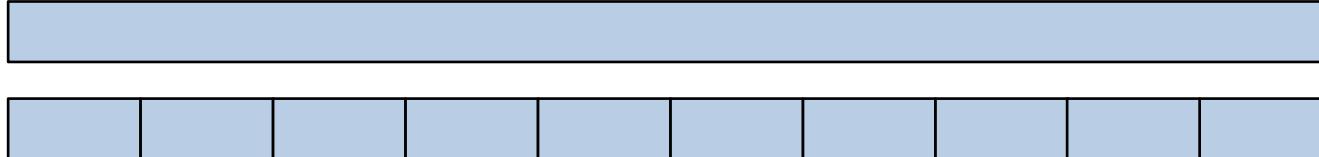


Micro Model Ensembles

- Train a **micro** model M_i on different **time slices** i

$$M_i = AD(X_i)$$

- 500 hours training data in total, intervals of 3-5 hours work well

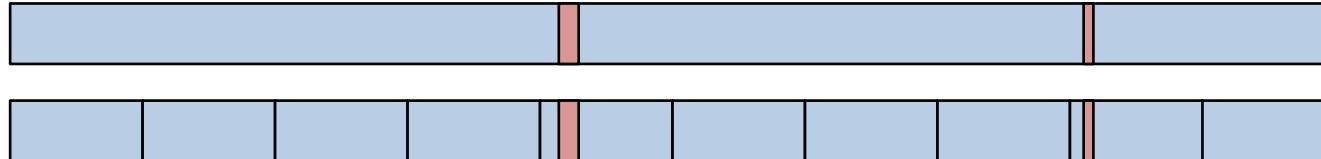


Micro Model Ensembles

- Train a **micro** model M_i on different **time slices** i

$$M_i = AD(X_i)$$

- 500 hours training data in total, intervals of 3-5 hours work well



Micro Model Ensembles

- Train a **micro** model M_i on different **time slices** i

$$M_i = AD(X_i)$$

- 500 hours training data in total, intervals of 3-5 hours work well
- **Sanitize** test package j by **combining output** of all micro models

$$\frac{1}{W} \sum_{i=1}^N w_i \times L_{i,j}$$

- **Set threshold** and divide data into benign and malicious
- **Train AD** on this **cleansed** data

Evaluation

- Use **two existing anomaly sensors** based on n-grams (but based on different learning algorithms)
- 300 hours for training, 100 for test

Sensor	www1		www		lists	
	FP(%)	TP(%)	FP(%)	TP(%)	FP(%)	TP(%)
A	0.07	0	0.01	0	0.04	0
A-S	0.04	20.20	0.29	17.14	0.05	18.51
A-SAN	0.10	100	0.34	100	0.10	100
P	0.84	0	6.02	40	64.14	64.19
P-SAN	6.64	76.76	10.43	61	2.40	86.54

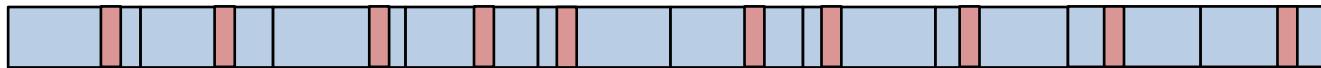
FP: false positive rate; TP: true positive rate

Evaluation of adaptive Attacks

- **Adaptive attacks**, evade by **introducing worm in every micromodel**
 - accuracy decreases to 30-40%

Evaluation of adaptive Attacks

- **Adaptive attacks**, evade by **introducing worm in every micromodel**
 - accuracy decreases to 30-40%



Evaluation of adaptive Attacks

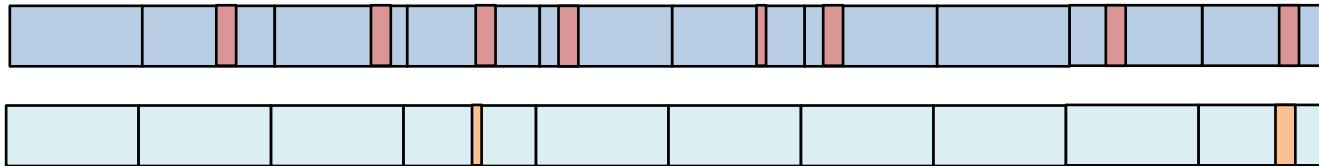
- **Adaptive attacks**, evade by **introducing worm in every micromodel**
 - accuracy decreases to 30-40%
- Solution: **share** models **across** different sites

$$M_{\text{cross}} = M_{\text{san}} - \bigcup \{M_{\text{abn}_i} \cap M_{\text{san}}\}$$

Evaluation of adaptive Attacks

- **Adaptive attacks**, evade by **introducing worm in every micromodel**
 - accuracy decreases to 30-40%
- Solution: **share** models **across** different sites

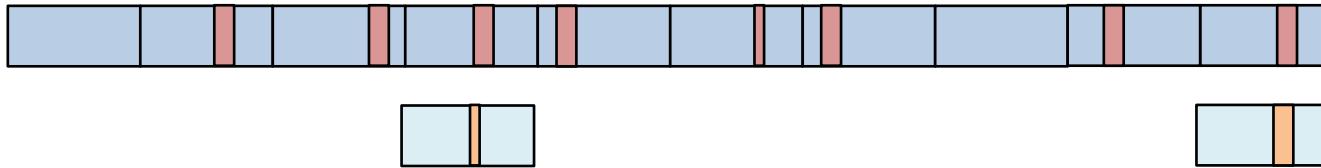
$$M_{\text{cross}} = M_{\text{san}} - \bigcup \{M_{\text{abn}_i} \cap M_{\text{san}}\}$$



Evaluation of adaptive Attacks

- **Adaptive attacks**, evade by **introducing worm in every micromodel**
 - accuracy decreases to 30-40%
- Solution: **share** models **across** different sites

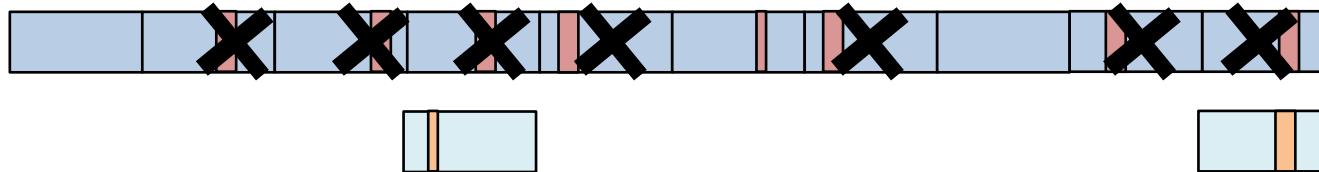
$$M_{\text{cross}} = M_{\text{san}} - \bigcup \{M_{\text{abn}_i} \cap M_{\text{san}}\}$$



Evaluation of adaptive Attacks

- **Adaptive attacks**, evade by **introducing worm in every micromodel**
 - accuracy decreases to 30-40%
- Solution: **share** models **across** different sites

$$M_{\text{cross}} = M_{\text{san}} - \bigcup \{M_{\text{abn}_i} \cap M_{\text{san}}\}$$



Evaluation of advanced Attacks

- **Adaptive attacks**, evade by **introducing worm in every micromodel**
 - accuracy decreases to 30-40%
- Solution: **share** models **across** different sites

Model	www1		www		lists	
	FP(%)	DR(%)	FP(%)	DR(%)	FP(%)	DR(%)
M_{pois}	0.10	44.94	0.27	51.78	0.25	47.53
M_{cross} (direct)	0.24	100	0.71	100	0.48	100
M_{cross} (indirect)	0.10	100	0.26	100	0.10	100

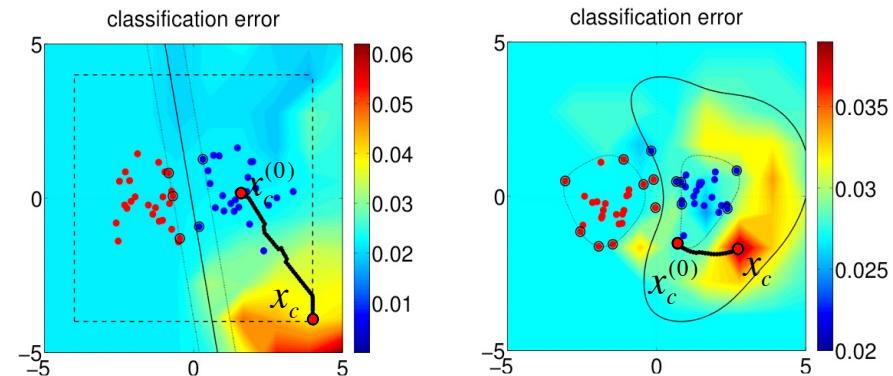
Take Aways



- Use knowledge from application
- Test adaptive attacks
- If possible, alter model to encompass attacks

Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - **Reject on negative impact**
- Robust Optimization
 - Bagging
 - PCA GRID
 - TRIM
- Beyond Defenses
 - Regularization and Poisoning
 - Adaptive Attacks
- Conclusion



Defense Summary – RONI

- Task:
 - Spam classification
- Algorithm:
 - Naive Bayes classifier
- Attack Assumptions:
 - Attacker's emails are always classified as spam, never as ham
 - Attacker can modify body of mail, but not header



Recap: Naive Bayes for Spam Classification

- Derived from **Bayes theorem**

$$p(H|D) = \frac{p(H)p(D|H)}{p(D)}$$



Recap: Naive Bayes for Spam Classification

- Derived from **Bayes theorem**
- In **practice**, instead of probabilities, use **word counts**:

$$P_S(w) = \frac{N_H N_S(w)}{N_H N_S(w) + N_S N_H(w)}$$

- Set **threshold** to define ham, spam, and unknown

Defense Strategy

- Reject On Negative Impact - RONI
- Training set of 20 mails, 5x50 in validation set.
- Train on training set plus a mail from validation set, **see if performance changes**



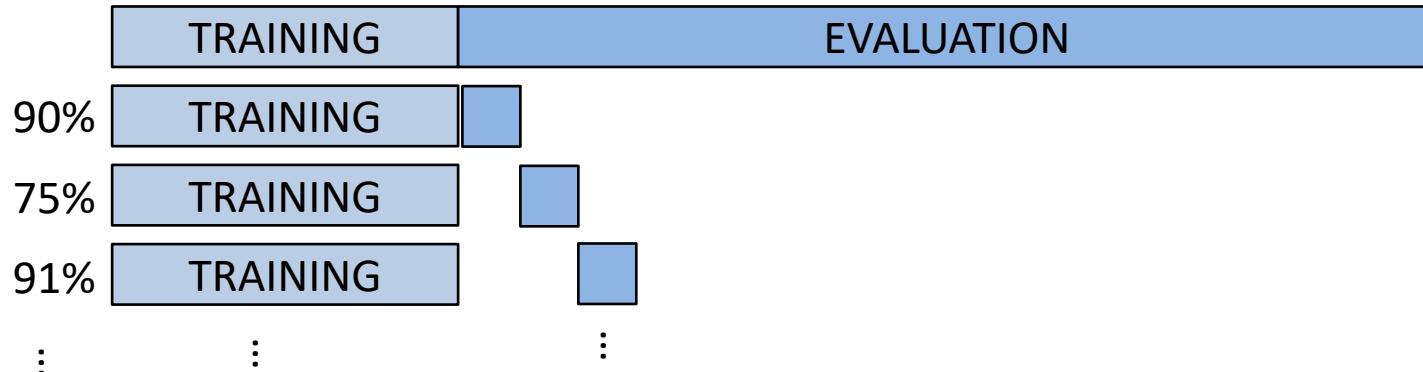
Defense Strategy

- Reject On Negative Impact - RONI
- Training set of 20 mails, 5x50 in validation set.
- Train on training set plus a mail from validation set, **see if performance changes**



Defense Strategy

- Reject On Negative Impact - RONI
- Training set of 20 mails, 5x50 in validation set.
- Train on training set plus a mail from validation set, **see if performance changes**

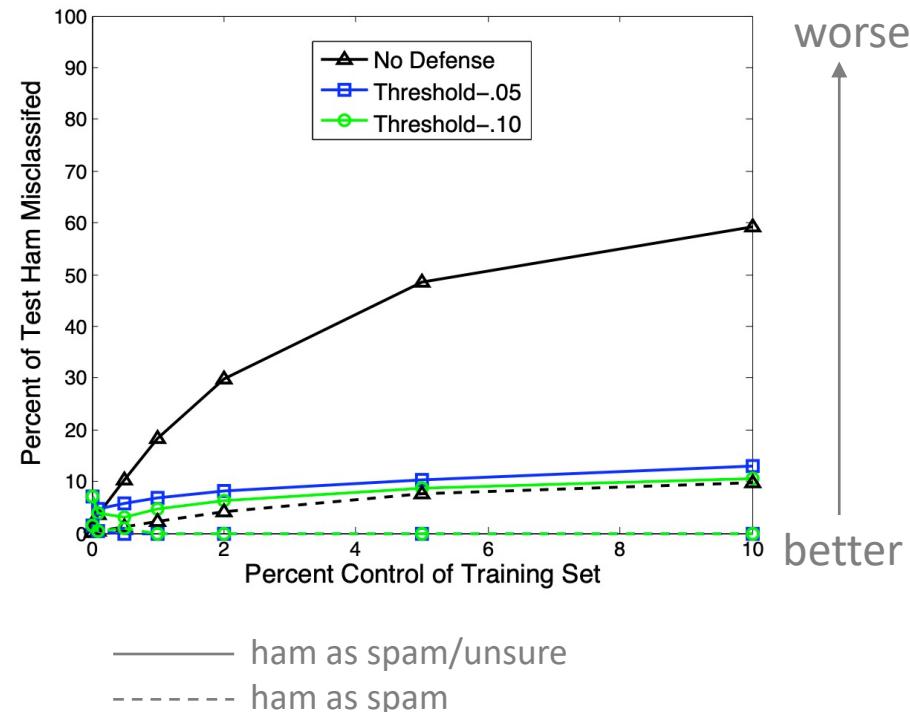


Defense Strategy

- Reject On Negative Impact - RONI
- Training set of 20 mails, 5x50 in validation set.
- Train on training set plus a mail from validation set, **see if performance changes**
- Attack message causes on average **~50% more decrease in true positives** as normal span
 - **Separable**
- Furthermore, thresholds can/should be adapted

Experiments

- 10, 000 inbox training set with 50% spam
- Problem: **correctly** identifies **ham**
- **fails** to identify **spam** (classified as unsure)



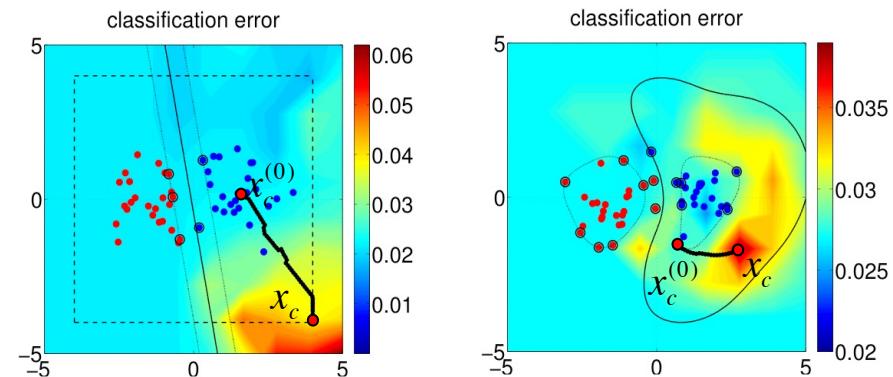
Take Aways



- Intuitive approach, but does not scale
- Add a whole dictionary in attack, impact has to be huge
- Dataset sanitization or robust learning?

Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - Reject on negative impact
- Robust Optimization
 - **Bagging**
 - PCA GRID
 - TRIM
- Beyond Defenses
 - Regularization and Poisoning
 - Adaptive Attacks
- Conclusion



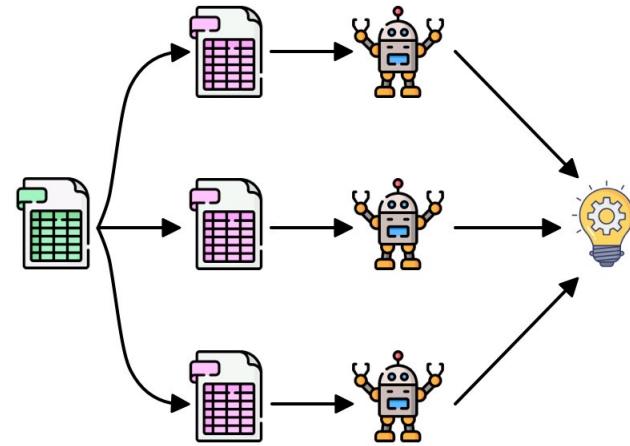
Defense Summary – Bagging

- Task:
 - Spam classification(, intrusion detection)
- Algorithm:
 - (weighted) Bagging
- Attack Assumptions:
 - Attacker controls subset of samples
 - Attacker changes features
- Further details:



Bagging and weighted bagging

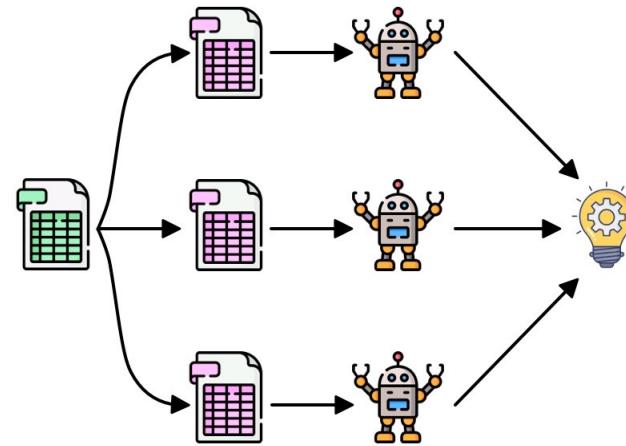
- Used for classifiers with high variance, e.g. large variation if data slightly different
- Bagging (e.g. bootstrap different datasets and then aggregate classifiers of each)



Bagging and weighted bagging

- Used for classifiers with high variance, e.g. large variation if data slightly different
- Bagging (e.g. bootstrap different datasets and then aggregate classifiers of each)

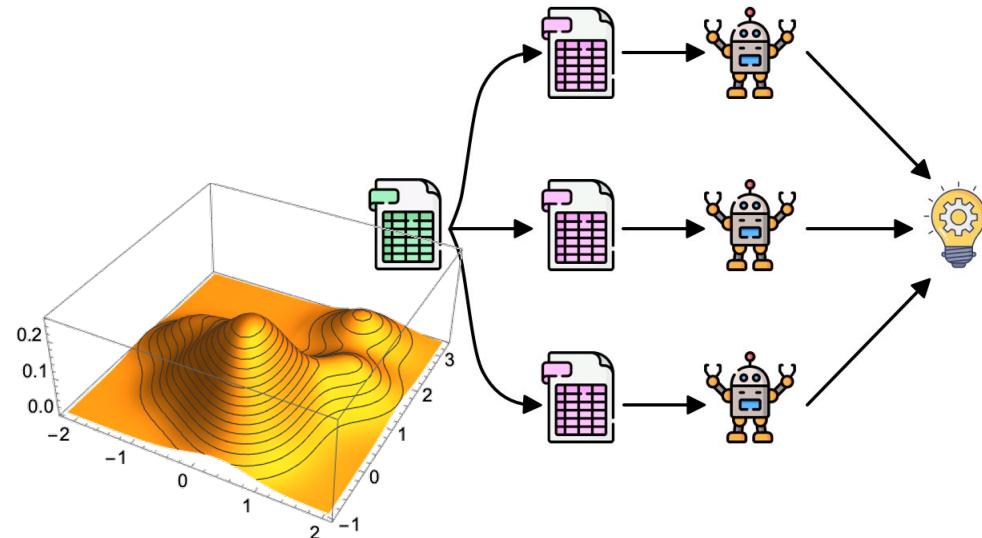
- Weighted bagging:
 - Density estimation on data
 - Weights are based on density:
high density -> high weight
low density -> low weight



Bagging and weighted bagging

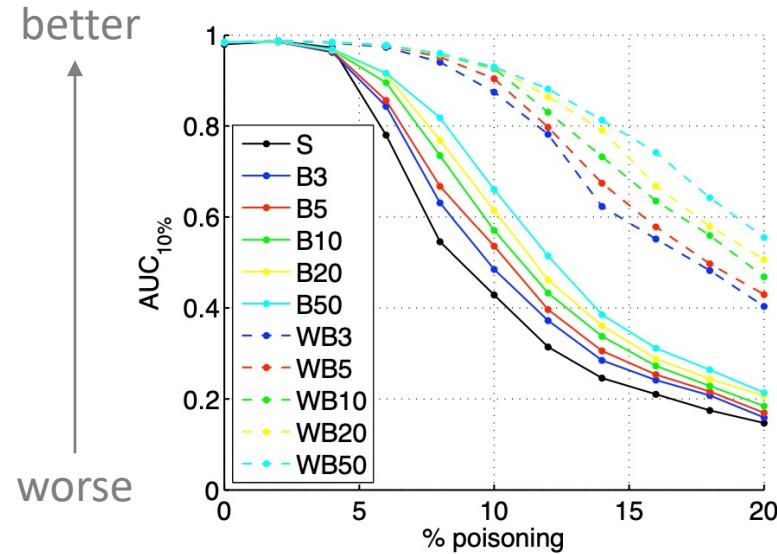
- Used for classifiers with high variance, e.g. large variation if data slightly different
- Bagging (e.g. bootstrap different datasets and then aggregate classifiers of each)

- Weighted bagging:
 - Density estimation on data
 - Weights are based on density:
high density \rightarrow high weight
low density \rightarrow low weight



Experiments on Spam

- TREC corpus, 25,220 legitimate and 50,199 spam emails
- 20,000 tokens, represented as binary



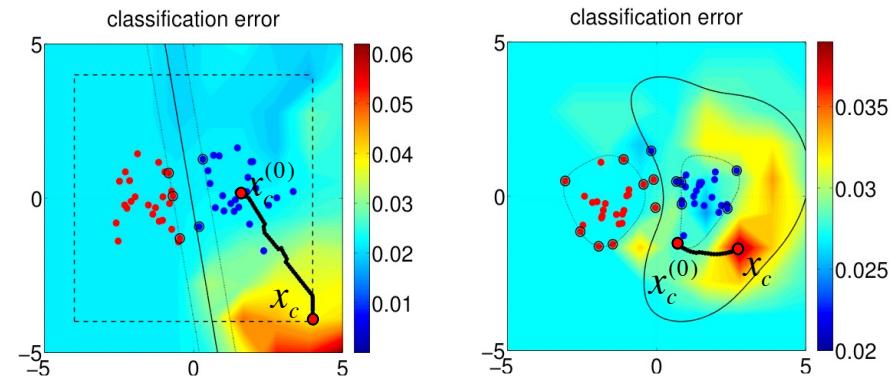
Take Aways



- RONI used random split, but here weighted bagging provides improvement
- Depends on kernel density estimation
- Depends on parameter which determines amount of classifiers

Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - Reject on negative impact
- Robust Optimization
 - Bagging
 - **PCA GRID**
 - TRIM
- Beyond Defenses
 - Regularization and Poisoning
 - Adaptive Attacks
- Conclusion



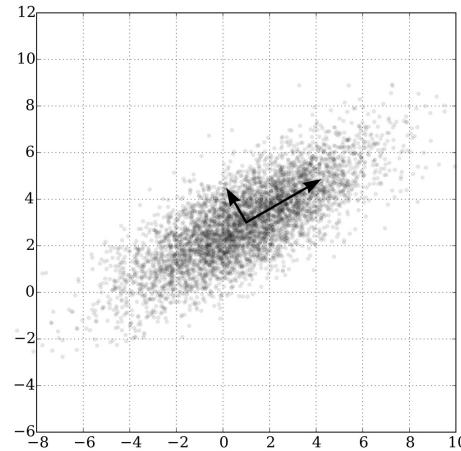
Defense Summary – PCA GRID

- Task:
 - Anomaly detection in backbone networks
- Algorithm:
 - PCA
- Attack Assumptions:
 - Attacker controls subset of samples
 - Attacker can only add traffic, not remove or delay
- Further details:
 - Study different attacks, White-box to black box, boiling frog



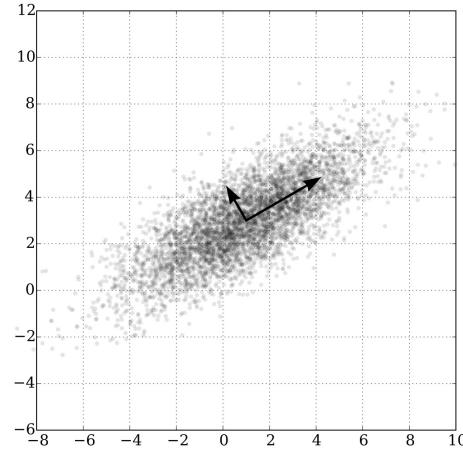
PCA and PCA-GRID

- PCA highly vulnerable to outliers:
- Based on mean, variance



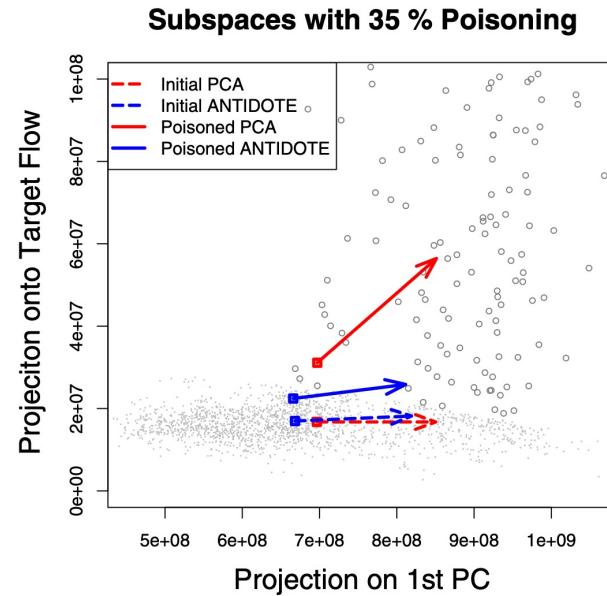
PCA and PCA-GRID

- PCA highly vulnerable to outliers
- Use instead Median Absolute Deviation (MAD)
- Use also median to center data



PCA and PCA-GRID

- PCA highly vulnerable to outliers
- Use instead Median Absolute Deviation (MAD)
- Use also median to center data

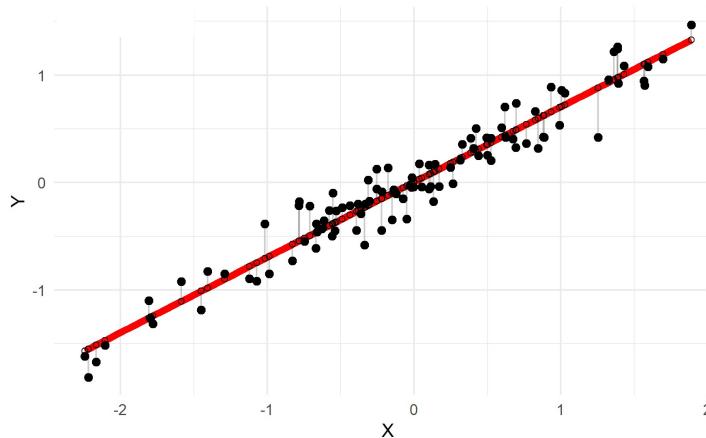
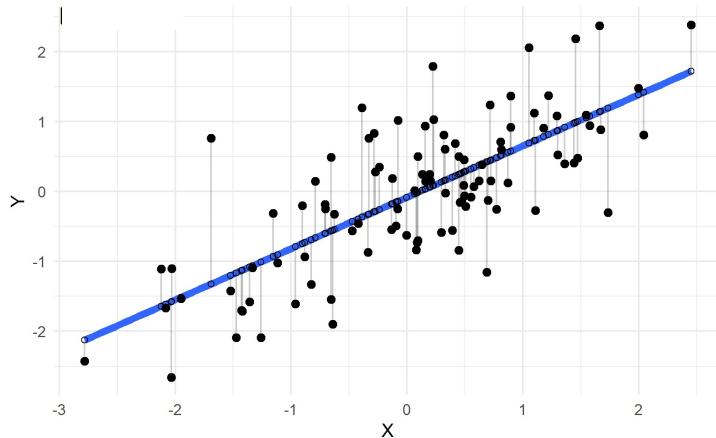


PCA residuals

- To detect malicious data, compose into normal and residual traffic

PCA residuals

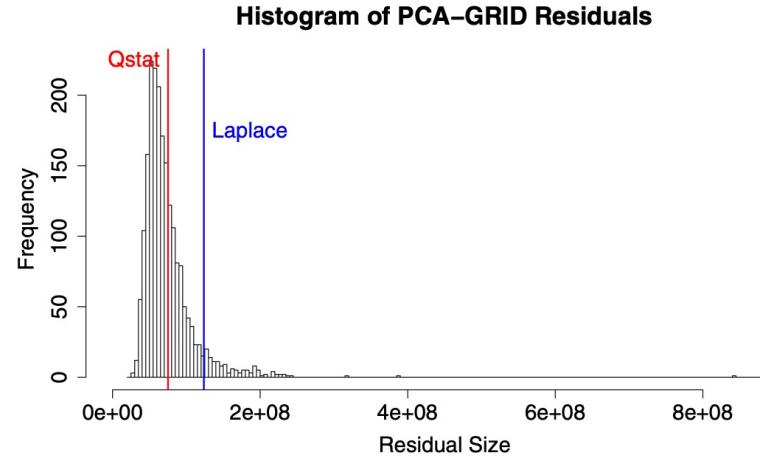
- To detect malicious data, compose into normal and residual traffic



- Small residuals: usual data, model fits well
- Large residuals: data does not fit well (or variance in data is large)

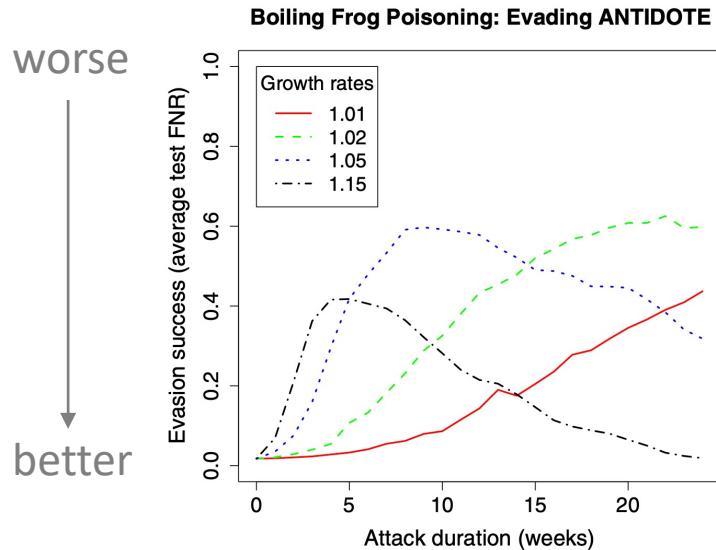
PCA residuals

- To detect malicious data, compose into normal and residual traffic
- Previous approaches cannot deal with heavy tailed data
- Chose Laplace cut-offs instead



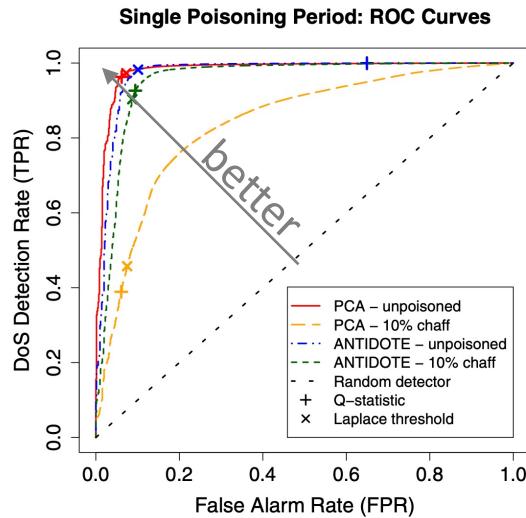
Experiments

- 6-month period from March 1 through September 10, 2004.
- Each week of data consists of 2016 measurements across 144 networks
- Evaluated on synthetic anomalies



Experiments

- 6-month period from March 1 through September 10, 2004.
- Each week of data consists of 2016 measurements across 144 networks
- Evaluated on synthetic anomalies



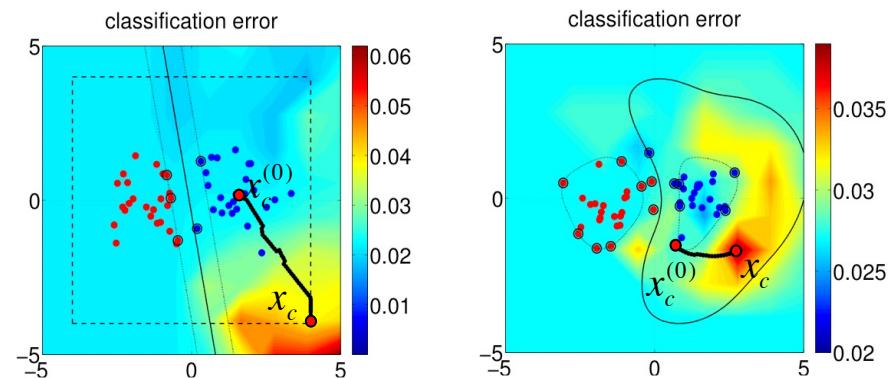
Take Aways



- Replace vulnerable components of algorithm
- Many problems from AML in general have already existing solutions in ML

Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - Reject on negative impact
- Robust Optimization
 - Bagging
 - PCA GRID
 - **TRIM**
- Beyond Defenses
 - Regularization and Poisoning
 - Adaptive Attacks
- Conclusion



Defense Summary – TRIM

- Task:
 - Loan, health care, and house price data
- Algorithm:
 - Linear regression
- Attack Assumptions:
 - Attacker controls subset of samples
 - Attacker changes features and labels



Approach

- Consider optimization problem:

$$\operatorname{argmin}_{w,b,I} L(w, b, I) = \frac{1}{|I|} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 + \lambda \Omega(\mathbf{w})$$
$$N = (1 + \alpha)n, \quad I \subset [1, \dots, N], \quad |I| = n$$

Approach

- Consider optimization problem:

$$\operatorname{argmin}_{w,b,I} L(w, b, I) = \frac{1}{|I|} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 + \lambda \Omega(\mathbf{w})$$
$$N = (1 + \alpha)n, \quad I \subset [1, \dots, N], \quad |I| = n$$

- Empirical risk minimization

Approach

- Consider optimization problem:

$$\operatorname{argmin}_{w,b,I} L(w, b, I) = \frac{1}{|I|} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 + \lambda \Omega(\mathbf{w})$$
$$N = (1 + \alpha)n, \quad I \subset [1, \dots, N], \quad |I| = n$$

- Empirical risk minimization + regularization

Approach

- Consider optimization problem:

$$\operatorname{argmin}_{w,b,I} L(w, b, I) = \frac{1}{|I|} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 + \lambda \Omega(\mathbf{w})$$
$$N = (1 + \alpha)n, \quad I \subset [1, \dots, N], \quad |I| = n$$

- Encompasses which points are used for optimization

Approach

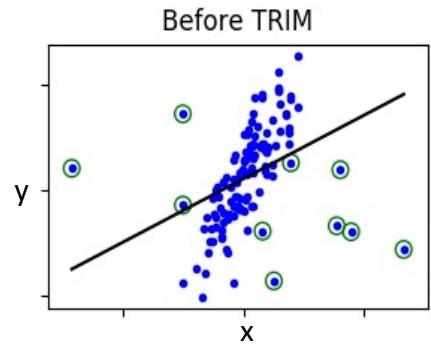
- Consider optimization problem:

$$\operatorname{argmin}_{w,b,I} L(w, b, I) = \frac{1}{|I|} \sum_{i \in I} (f(\mathbf{x}_i) - y_i)^2 + \lambda \Omega(\mathbf{w})$$
$$N = (1 + \alpha)n, \quad I \subset [1, \dots, N], \quad |I| = n$$

- Difficult to solve, so iterative solution
- Empirically shown that converges. Consider that no assumptions on how attack is generated (just on amount of poisoned samples)

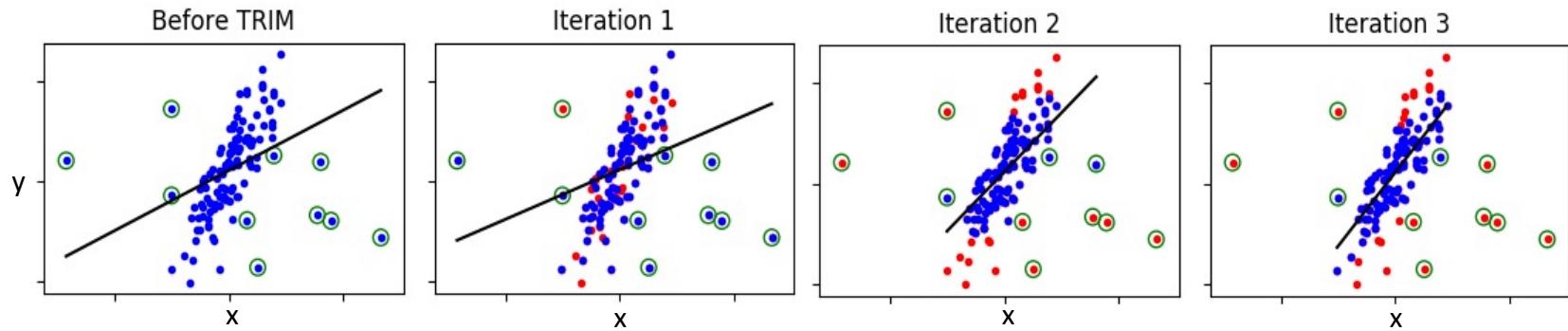
TBD

- High Level intuition of algorithm



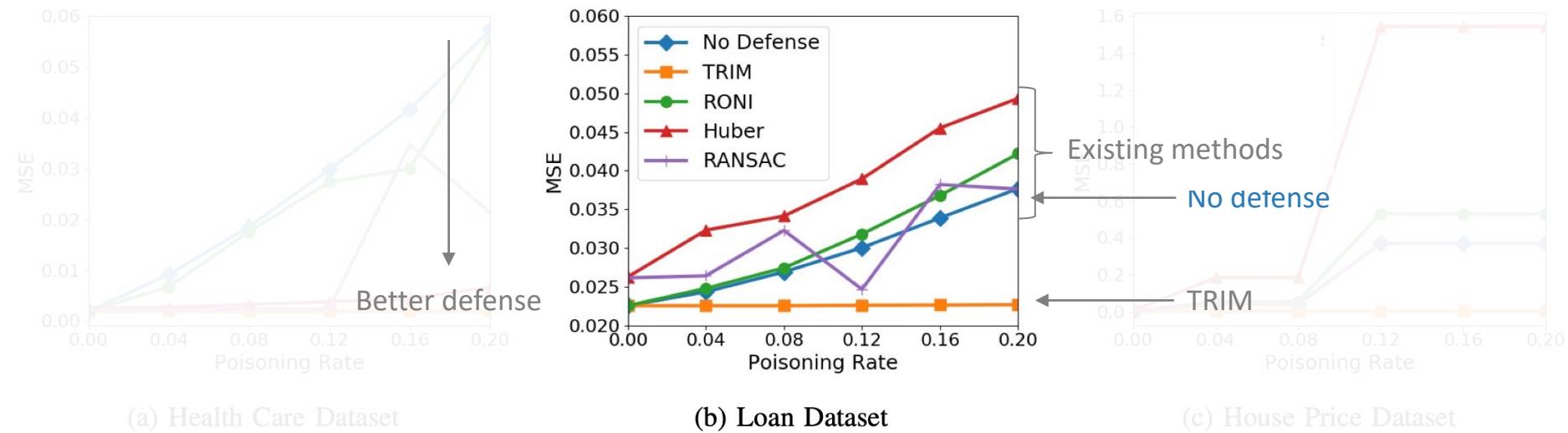
TBD

- High Level intuition of algorithm



Evaluation

- 3 datasets, health care, loan, and house pricing



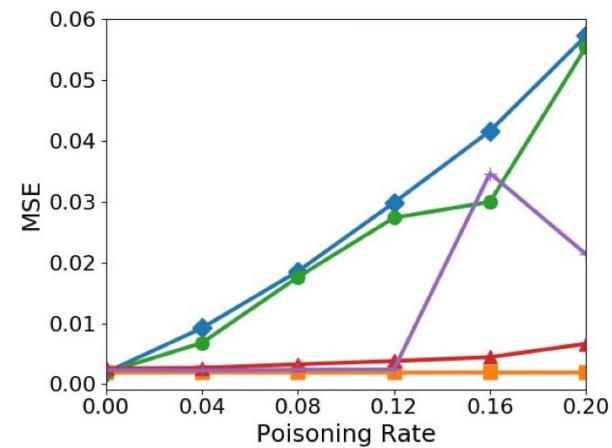
(a) Health Care Dataset

(b) Loan Dataset

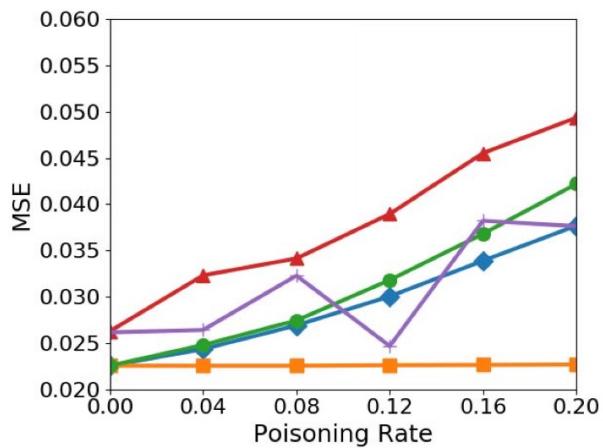
(c) House Price Dataset

Evaluation

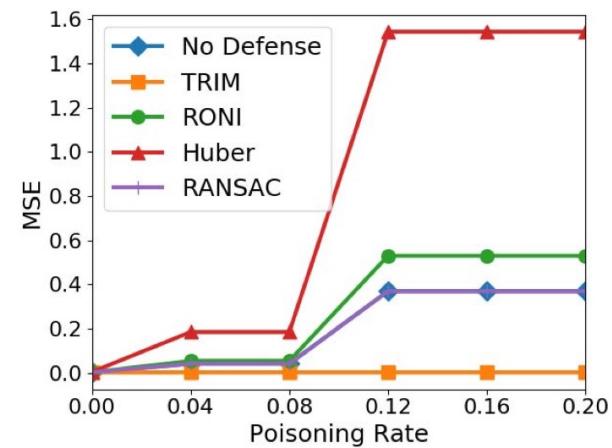
- 3 datasets, health care, loan, and house pricing



(a) Health Care Dataset



(b) Loan Dataset



(c) House Price Dataset

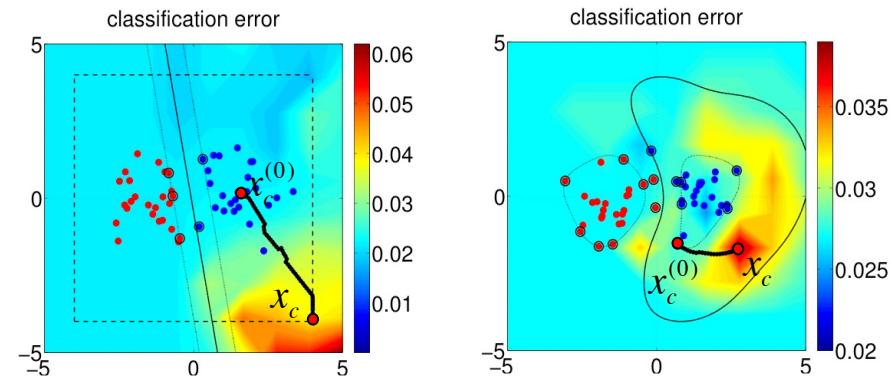
Take Aways



- Outlier detection is directly incorporated into optimization problem
- Compare to related work
- Amount of introduced poisoned samples must be limited

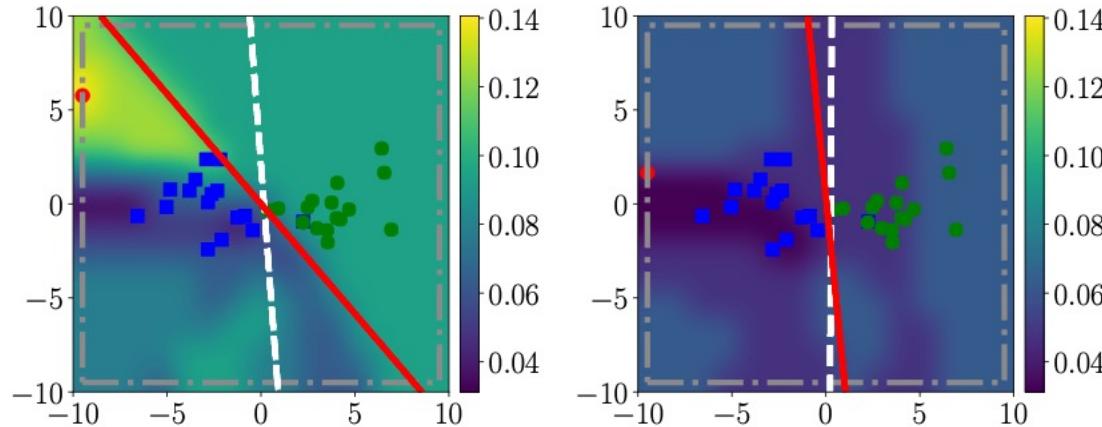
Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - Reject on negative impact
- Robust Optimization
 - Bagging
 - PCA GRID
 - TRIM
- Beyond Defenses
 - **Regularization and Poisoning**
 - Adaptive Attacks
- Conclusion

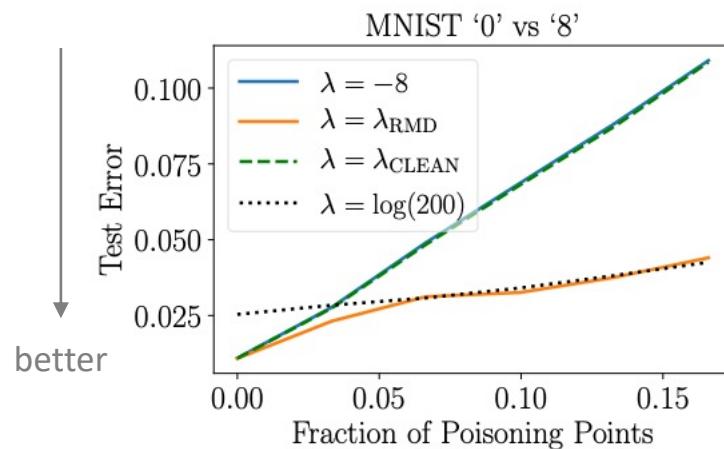


Regularization and Poisoning

- We have been talking about limiting the influence of a single point on the algorithm
 - Existing work on that, called regularization

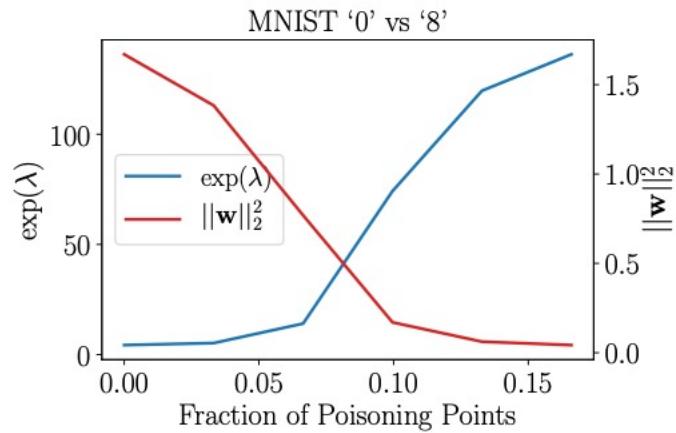


Learning (poisoned?) hyperparameters



- Idea: run hyperparameter search to find good hyperparameters for regularization, also monitor weights
 - When param is learned, error is lower
- Effect of poisoning can be alleviated with proper hyperparameters

Learning (poisoned?) hyperparameters II



- More poisoning points imply larger learned regularization hyperparameter
- Consequently, size of the weights get smaller

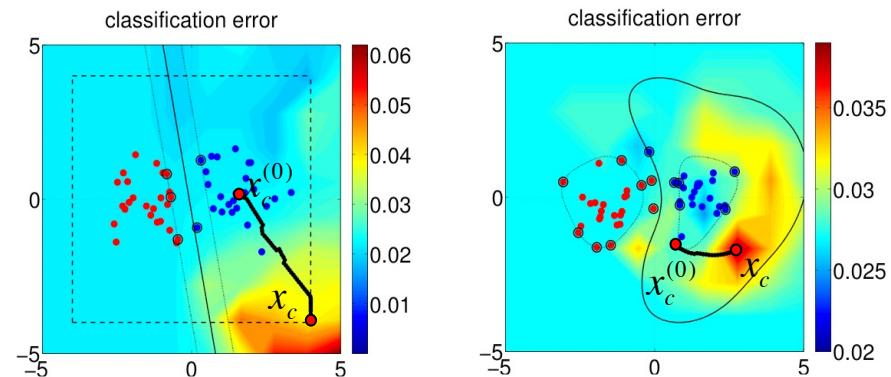
Take Aways



- Regularization affects how the classifier learns the data, and thus the poison
- relationship between complexity and poison vulnerability:
 - a flexible classifier learns poisons faster
 - a less flexible classifier learns poisons slower
- For many problems from security, ML solutions exist already!

Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - Reject on negative impact
- Robust Optimization
 - Bagging
 - PCA GRID
 - TRIM
- Beyond Defenses
 - Regularization and Poisoning
 - **Adaptive Attacks**
- Conclusion



Evaluating Defenses

- Important to conduct ablation studies on parameters
- State assumptions / threat model clearly

Evaluating Defenses

- Important to conduct ablation studies on parameters
- State assumptions / threat model clearly
- Most important, consider adaptive attacks...
- ... knowing the defense, how would an attack target the system?

Setting for an adaptive Attack

- Data sanitization-based defenses:
 - Outlier detector that rejects anything outside a radius (L2 defense)
 - exclude points with large loss (TRIM, not tested)
 - ... more in the paper
- Assumptions:
 - defense is applied without human intervention
 - Attacker can add 3%, remove 5%
- The paper introduces several attacks, we focus on one.

Basic Minimax Attack

- To approximate test loss, training data is used
- Essentially solve a saddle point problem

$$\underset{\mathcal{D}_p \subseteq \mathcal{F}_\beta}{\text{maximize}} \quad L(\theta; \mathcal{D}_c \cup \mathcal{D}_p)$$

$$\text{where} \quad \hat{\theta} \stackrel{\text{def}}{=} \underset{\theta}{\operatorname{argmin}} L(\theta; \mathcal{D}_c \cup \mathcal{D}_p)$$

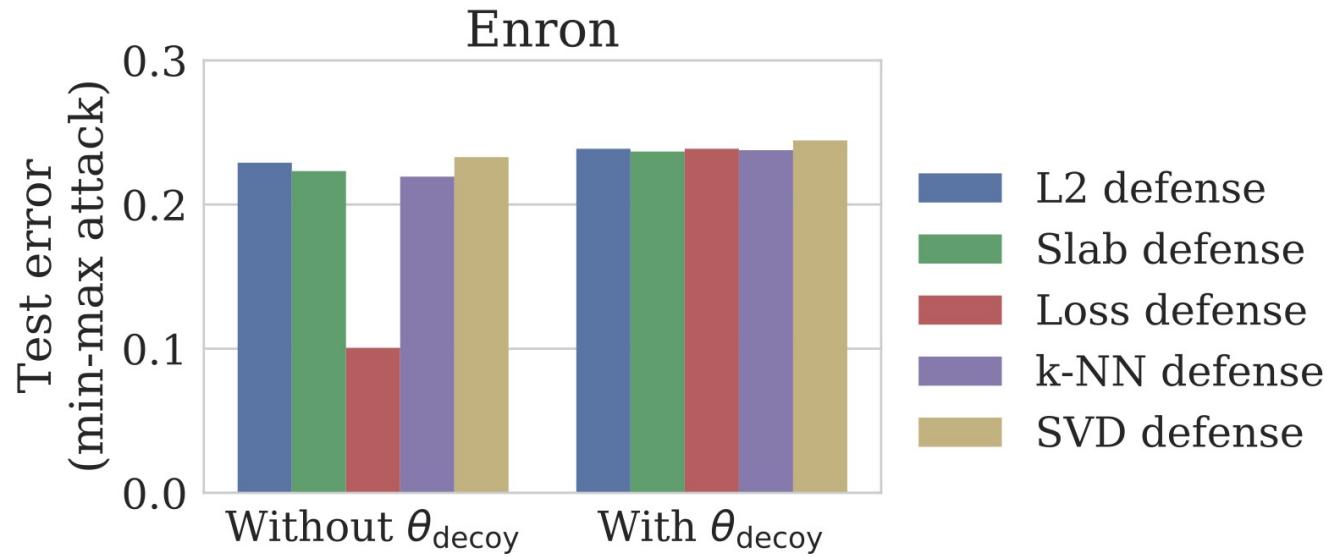
- Convex if the loss is convex

Adaptive Minimax

- Issues with old attack:
 - Approximated test loss might be wrong because model does not fit data well
 - High loss points are detectable by the loss defense
- Solution: add constraint where decoy parameters' loss on the poisoning data have low loss
 - Decoy paramters have high test loss, real params are driven towards them
 - Upper bounds loss value of poisoning points

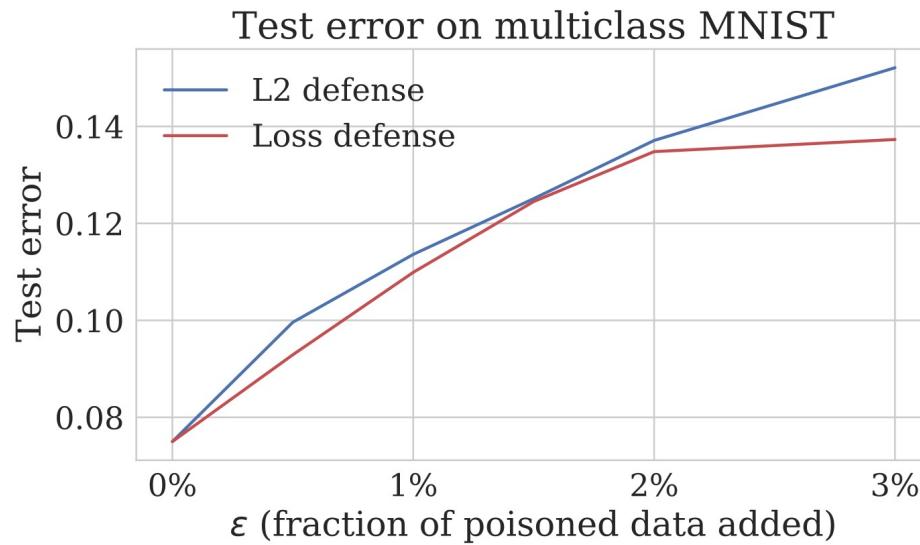
Empirical results

- Comparing new and adapted attack on the Enron dataset



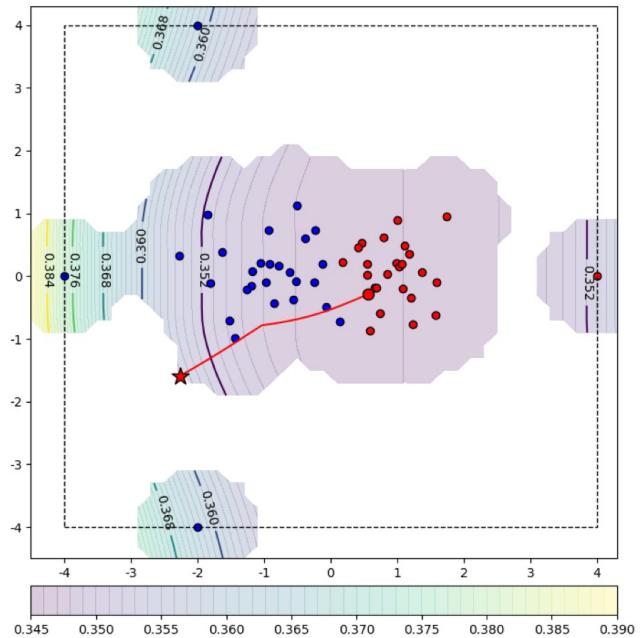
Empirical results

- Increasing the amount of poisoning points on defenses

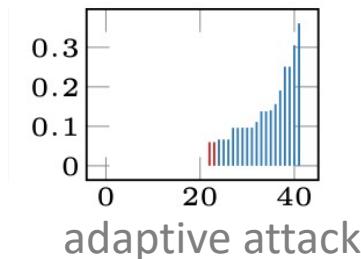
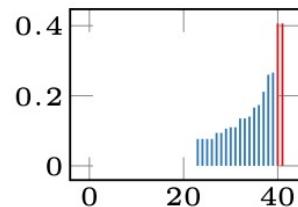
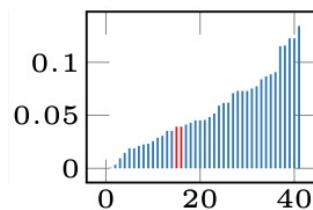
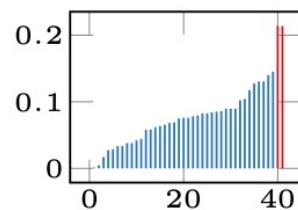
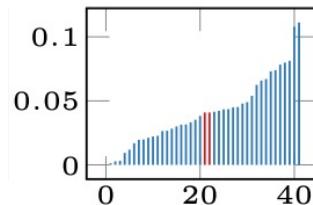
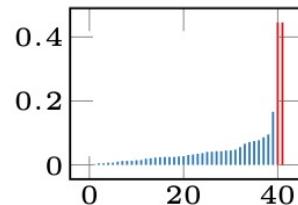


Perturbation Constraints

- Benign attacks lead to points that are far away from the data
- Easy to detect
- Solution: Attacker distance threshold limits how far attacks can be maximally away from the original points



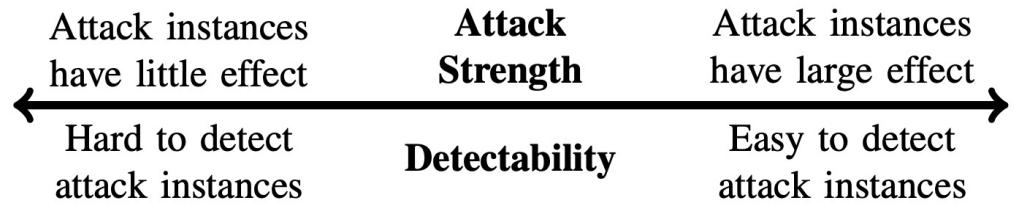
Perturbation Constraints in Practice



- Score based on DT trained to separate outlier points
- Three different datasets
- Benign data blue, red outlier scores

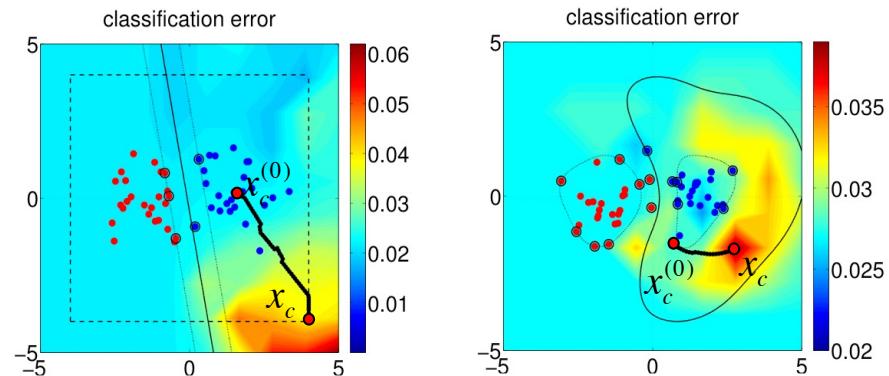
Take Aways

- Evaluation on adaptive attack ABSOLUTELY necessary
- High level insight:



Countering Poisoning Attacks - Outline

- Data Sanitization Defenses
 - Micro models
 - Reject on negative impact
- Robust Optimization
 - Bagging
 - PCA GRID
 - TRIM
- Beyond Defenses
 - Regularization and Poisoning
 - Adaptive Attacks
- Conclusion



Conclusion

- In contrast to evasion, hope for defenses
- Roughly two groups, data sanitization and robust optimization
- Few strongly changed (and easily detectable) points or more, less changed (and thus hard to detect) points.
- A flexible classifier learns poisoning points fast, whereas a less flexible classifier will lead to worse poisoning success



Kathrin Grosse

Kathrin.grosse@unica.it

 Kathrin Grosse

Thanks!



What is the rule? The rule is protect yourself at all times
(from the movie "Million dollar baby", 2004)