

AI 539: TRUSTWORTHY ML **(CERTIFIED) DEFENSES AGAINST POISONING ATTACKS**

Sanghyun Hong

sanghyun.hong@oregonstate.edu



Oregon State
University

SAIL

Secure AI Systems Lab

Traditionally, computer security seeks to ensure a system's integrity against attackers by creating clear boundaries between the system and the outside world (Bishop, 2002). In machine learning, however, the most critical ingredient of all—the training data—comes directly from the outside world.

– Steinhardt, Koh, and Liang, NeurIPS'17

DEFENSES AGAINST DATA POISONING ATTACKS

- Existing defenses
 - RONI (Reject on Negative Impact)
 - TRIM
 - tRONI¹
 - ... (many more)
- Problem:
 - Existing defenses empirically works
 - How can we provide “provable” defense guarantee against poisoning attacks?

DEFENSES AGAINST DATA POISONING ATTACKS

- What we “provably” guarantee?
 - A model’s loss over the test-set (or a subset of it) is less than a specific value
 - The above is valid when the # of poisons in the training data are less than a specific value
- What are the types of “provable” defenses?
 - Pre-training defense: data sanitization
 - Training-time defense: novel training algorithms

“PROVABLE” DATA SANITIZATION DEFENSE

CERTIFIED DEFENSES FOR DATA POISONING ATTACKS, STEINHARDT ET AL., NEURIPS 2017

THREAT MODEL

- Setup [binary classification task!]
 - **Data:** $x \in X$ (ex. R^d), $y \in Y = \{-1, +1\}$
 - **Clean train-set:** D_c of size n / **Test-set:** S
 - **Loss function:** $l(\theta; x, y) = \max(0, 1 - y\langle\theta, x\rangle)$
 - **Test-loss:** $L(\theta) = \mathbb{E}_{(x,y) \sim S}[l(\theta; x, y)]$

- Data sanitization defenses

- **Goal:** Examine $D_c \cup D_p$ and remove poisons (e.g., outliers)

$$\hat{\theta} \stackrel{\text{def}}{=} \underset{\theta \in \Theta}{\operatorname{argmin}} L(\theta; (\mathcal{D}_c \cup \mathcal{D}_p) \underbrace{\cap \mathcal{F}}_{\text{A filtered dataset}}), \quad \text{where } L(\theta; S) \stackrel{\text{def}}{=} \sum_{(x,y) \in S} \ell(\theta; x, y)$$

- **Methods:**

- *Fixed* (oracle) defense: when we know the true distribution of data (unrealistic)
 - *Data-dependent* defense: when we don't know the true distribution (real-world!)

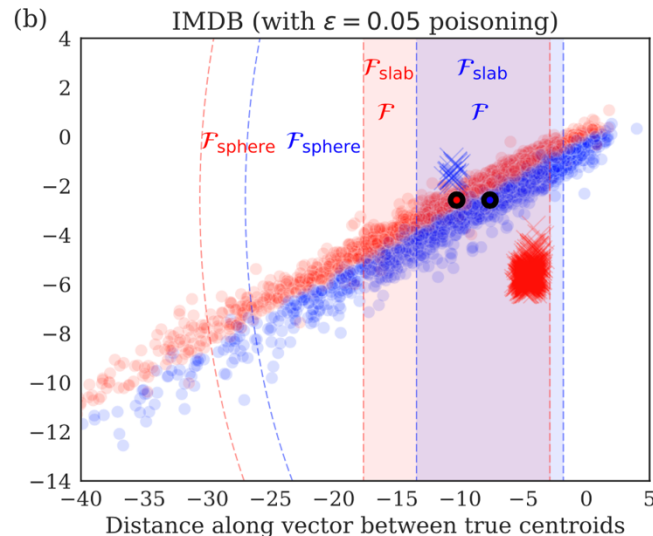
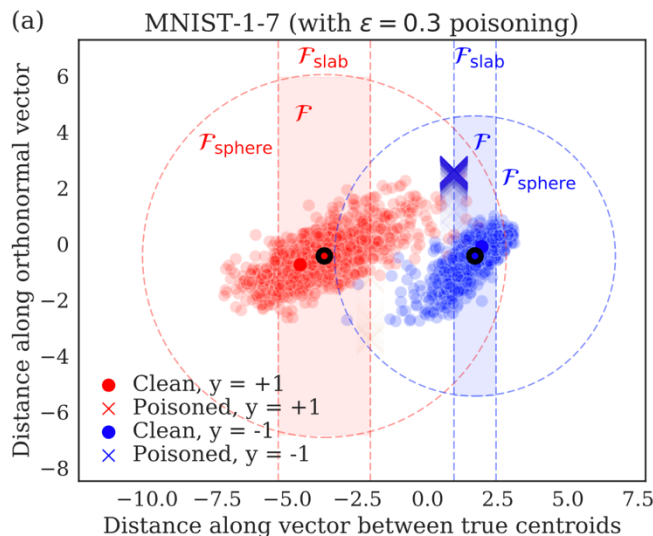
EXAMPLE DATA SANITIZATION DEFENSES

- Data sanitization defenses

- **Goal:** Examine $D_c \cup D_p$ and remove poisons (e.g., outliers)

- **Example defenses:**

- *sphere* defense: removes points outside a spherical radius
 - *slab* defense: first project points onto the line btw. the centroids and then remove



THE WORST-CASE TEST LOSS UNDER DATA POISONING

- Upper-bound [refer to the paper for its derivation]

$$\max_{D_p} L(\hat{\theta}) \leq \max_{D_p \subseteq F} \min_{\theta \in \Theta} \frac{1}{n} L(\theta; D_c \cup D_p) \stackrel{\text{def}}{=} \mathbf{M}$$

- **M**: the minimax loss
- **It means**: the attack is bounded to a scenario where **all poisons are alive**!

THE WORST-CASE TEST LOSS WITH A DEFENSE F

- Upper-bound [refer to the paper for its derivation]

$$\max_{D_p} L(\hat{\theta}) \leq \max_{D_p \subseteq F} \min_{\theta \in \Theta} \frac{1}{n} L(\theta; D_c \cup (D_p \cap F)) \stackrel{\text{def}}{=} \mathbf{M}$$

- **M**: the minimax loss
 - **It means**: the attack is bounded to a scenario where **all poisons are alive** under F !
-
- Two defense scenarios
 - **Fixed defense**: when we know the true distribution of data
 - **Data-dependent defense**: when we don't know the true distribution of data

THE WORST-CASE TEST LOSS WITH A FIXED DEFENSE

- Upper-bound [refer to the paper for its derivation]

$$\max_{D_p} L(\hat{\theta}) \leq \max_{D_p \subseteq F} \min_{\theta \in \Theta} \frac{1}{n} L(\theta; D_c \cup (D_p \cap F)) \stackrel{\text{def}}{=} \mathbf{M}$$

- **M**: the minimax loss
 - **It means**: the attack is bounded to a scenario where **all poisons are alive** under **F**!
-
- Two defense scenarios
 - **Fixed defense**: we can **fix F** regardless of poisoning samples
 - **Data-dependent defense**: when we don't know the true distribution of data

HOW DO WE COMPUTE THE UPPER-BOUND FOR A FIXED DEFENSE?

- *Fixed* defense scenario
 - To compute the upper-bound, you iteratively craft poisons and train models on them

Algorithm 1 Online learning algorithm for generating an upper bound and candidate attack.

Input: clean data \mathcal{D}_c of size n , feasible set \mathcal{F} , radius ρ , poisoned fraction ϵ , step size η .

Initialize $z^{(0)} \leftarrow 0$, $\lambda^{(0)} \leftarrow \frac{1}{\eta}$, $\theta^{(0)} \leftarrow 0$, $U^* \leftarrow \infty$.

for $t = 1, \dots, \epsilon n$ **do**

 Compute $(x^{(t)}, y^{(t)}) = \operatorname{argmax}_{(x,y) \in \mathcal{F}} \ell(\theta^{(t-1)}; x, y)$.

$U^* \leftarrow \min(U^*, \frac{1}{n}L(\theta^{(t-1)}; \mathcal{D}_c) + \epsilon \ell(\theta^{(t-1)}; x^{(t)}, y^{(t)}))$.

$g^{(t)} \leftarrow \frac{1}{n} \nabla L(\theta^{(t-1)}; \mathcal{D}_c) + \epsilon \nabla \ell(\theta^{(t-1)}; x^{(t)}, y^{(t)})$.

 Update: $z^{(t)} \leftarrow z^{(t-1)} - g^{(t)}$, $\lambda^{(t)} \leftarrow \max(\lambda^{(t-1)}, \frac{\|z^{(t)}\|_2}{\rho})$, $\theta^{(t)} \leftarrow \frac{z^{(t)}}{\lambda^{(t)}}$.

} Iteratively craft poisons
to fool the t -th classifier

end for

Output: upper bound U^* and candidate attack $\mathcal{D}_p = \{(x^{(t)}, y^{(t)})\}_{t=1}^{\epsilon n}$.

- **Proposition:** $U^* - \frac{1}{n}L(\tilde{\theta}; \mathcal{D}_c \cup \mathcal{D}_p) \leq \frac{\operatorname{Regret}(\epsilon n)}{\epsilon n}$

Any poisoning that minimizes the avg. Regret will be close to the optimal

THE WORST-CASE TEST LOSS WITH A DATA-DEPENDENT DEFENSE

- Upper-bound [refer to the paper for its derivation]

$$\max_{D_p} L(\hat{\theta}) \leq \max_{D_p \subseteq F} \min_{\theta \in \Theta} \frac{1}{n} L(\theta; D_c \cup (D_p \cap F)) \stackrel{\text{def}}{=} \mathbf{M}$$

- **M**: the minimax loss
 - **It means**: the attack is bounded to a scenario where **all poisons are alive** under **F**!
-
- Two defense scenarios
 - **Fixed defense**: we can **fix F** regardless of poisoning samples
 - **Data-dependent defense**: we **cannot fix F** (and hence can be influenced by the attacker)

HOW DO WE COMPUTE THE UPPER-BOUND FOR A DATA-DEP. DEFENSE?

- *Data-dependent* defense scenario
 - ex. In Slab defense, one can use the **empirical mean** instead of the true mean

Algorithm 1 Online learning algorithm for generating an upper bound and candidate attack.

Input: clean data \mathcal{D}_c of size n , feasible set \mathcal{F} , radius ρ , poisoned fraction ϵ , step size η .

Initialize $z^{(0)} \leftarrow 0$, $\lambda^{(0)} \leftarrow \frac{1}{\eta}$, $\theta^{(0)} \leftarrow 0$, $U^* \leftarrow \infty$.

for $t = 1, \dots, \epsilon n$ **do**

 Compute $(x^{(t)}, y^{(t)}) = \operatorname{argmax}_{(x,y) \in \mathcal{F}} \ell(\theta^{(t-1)}; x, y)$.

$U^* \leftarrow \min(U^*, \frac{1}{n}L(\theta^{(t-1)}; \mathcal{D}_c) + \epsilon \ell(\theta^{(t-1)}; x^{(t)}, y^{(t)}))$.

$g^{(t)} \leftarrow \frac{1}{n} \nabla L(\theta^{(t-1)}; \mathcal{D}_c) + \epsilon \nabla \ell(\theta^{(t-1)}; x^{(t)}, y^{(t)})$.

 Update: $z^{(t)} \leftarrow z^{(t-1)} - g^{(t)}$, $\lambda^{(t)} \leftarrow \max(\lambda^{(t-1)}, \frac{\|z^{(t)}\|_2}{\rho})$, $\theta^{(t)} \leftarrow \frac{z^{(t)}}{\lambda^{(t)}}$.

end for

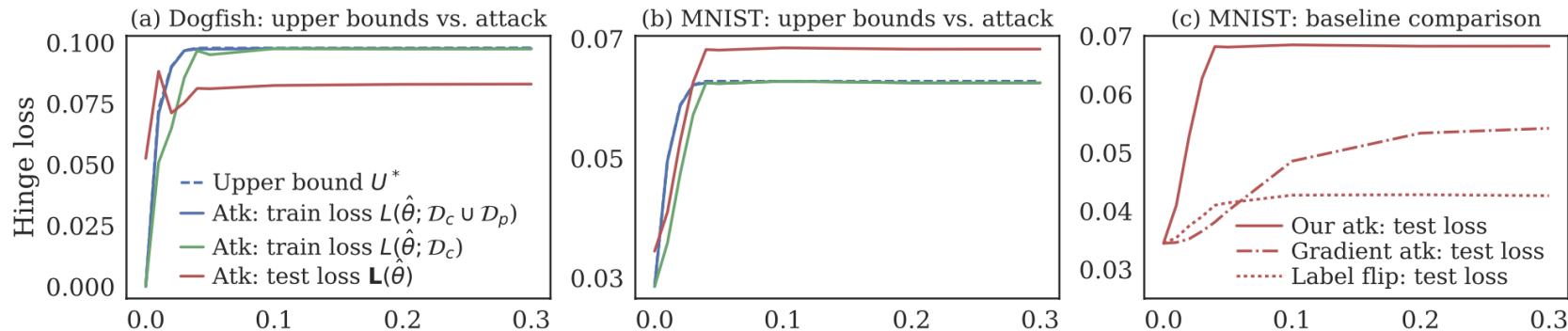
Output: upper bound U^* and candidate attack $\mathcal{D}_p = \{(x^{(t)}, y^{(t)})\}_{t=1}^{\epsilon n}$.

- **Proposition:** $\tilde{U}(\theta) \stackrel{\text{def}}{=} \frac{1}{n}L(\theta; \mathcal{D}_c) + \epsilon \max_{\operatorname{supp}(\pi_p) \subseteq \mathcal{F}(\pi_p)} \mathbf{E}_{\pi_p}[\ell(\theta; x, y)]$

Any poisoning that minimizes the avg. Regret will be **close to the optimal**
Here we estimate the Regret over any probability distribution π_p

EVALUATIONS: UNDER A FIXED DEFENSE F

- On DogFish and MNIST-1/7



- **Notations:**

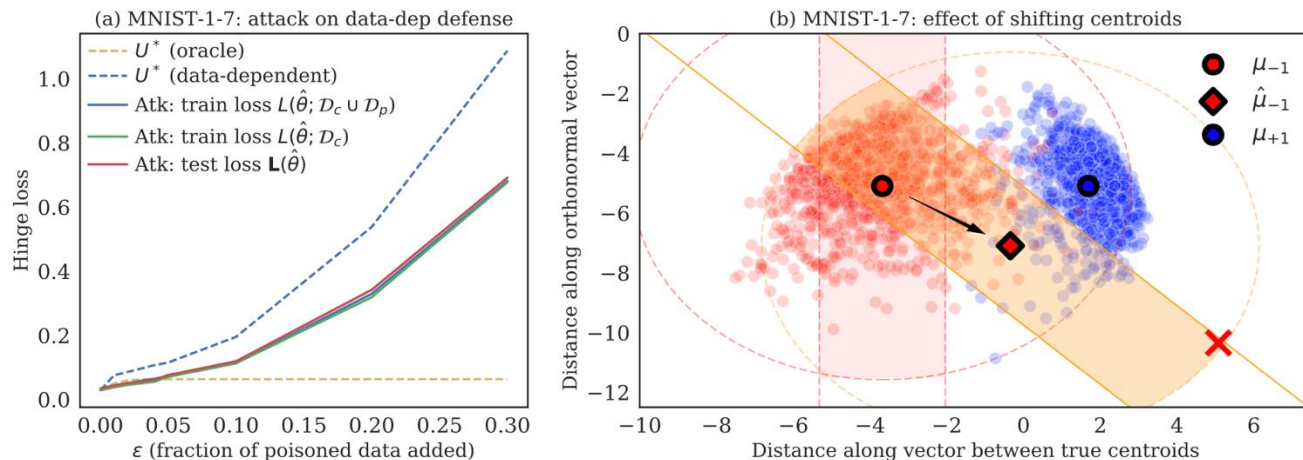
- (solid blue) the candidate attack | (dashed blue) the worst-case train loss (Prep.)

- **Takeaways:**

- (a), (b), (c): the fixed defense is strong (the loss < 0.1...)
- (a) and (b): the upper bound is *tight*
- (c): the upper bound is tighter than what existing attacks can inflict

EVALUATIONS: UNDER A DATA-DEPENDENT DEFENSE F

- On MNIST-1/7 in 2-class SVMs



- (a): data-dependent defenses are much weaker (the bound increases exponentially...)
- (a): the upper-bound is still *tight*
- (b): in data-dependent defenses, the F is affected by the poisons

“PROVABLE” TRAINING-TIME DEFENSE

DATA POISONING AGAINST DIFFERENTIALLY-PRIVATE LEARNERS: ATTACKS AND DEFENSES, MA ET AL., IJCAI 2019

TRAINING-TIME DEFENSES

- Desiderata
 - A defense wants to reduce a model's sensitivity to the training data alterations
 - More precisely
 - D is a training set drawn from the data distribution
 - \tilde{D} is a compromised training set, by an adversary
 - f is a model, and f_D and $f_{\tilde{D}}$ are the models trained on D and \tilde{D}
 - f_D and $f_{\tilde{D}}$ behave similarly (or the same) on the test-set

DIFFERENTIAL PRIVACY

- ϵ -Differential Privacy

- A randomized algorithm $M: D \rightarrow R$ with domain D and a range R satisfies ϵ -differential privacy if for any two adjacent inputs $d, d' \in D$ and any subset of outputs $S \subset R$ it holds

$$\Pr[\mathcal{M}(d) \in S] \leq e^\epsilon \Pr[\mathcal{M}(d') \in S]$$

- (ϵ, δ) -Differential Privacy

$$\Pr[\mathcal{M}(d) \in S] \leq e^\epsilon \Pr[\mathcal{M}(d') \in S] + \delta$$

- δ : Represent some catastrophic failure cases [[Link](#), [Link](#)]
- $\delta < 1/|d|$, where $|d|$ is the number of samples in a database

DIFFERENTIAL PRIVACY

- (ϵ, δ) -Differential Privacy [Conceptually]

$$\Pr[\mathcal{M}(d) \in S] \leq e^\epsilon \Pr[\mathcal{M}(d') \in S] + \delta$$

- You have two databases d, d' differ by one item
- You make the same query M to each and have results $M(d)$ and $M(d')$
- You ensure the distinguishability between the two under a measure ϵ
 - ϵ is large: those two are distinguishable, less private
 - ϵ is small: the two outputs are similar, more private
- You also ensure the catastrophic failure probability under δ

DIFFERENTIAL PRIVACY

- (ϵ, δ) -Differential Privacy
 - Implementation: Gaussian mechanism
 - **Formally:**
 - Suppose properties $q = (q_1, \dots, q_k)$
 - Gaussian mechanism M_{q, σ^2} takes
 - >> x as input (or gradients as input)
 - >> releases $\hat{q} = (\hat{q}_1, \dots, \hat{q}_k)$
 - where each \hat{q}_i is independent sample from $N(q_i(x), \sigma^2)$,
 - for an appropriate variance σ^2
 - **Easy-way:**
 - Add Gaussian noise with a variance σ^2 to
 - >> the output \hat{q} (output perturbation)
 - >> the gradients (object perturbation)
 - such that the output satisfies ϵ -differential privacy guarantee

TRAINING-TIME DEFENSES: THREAT MODEL

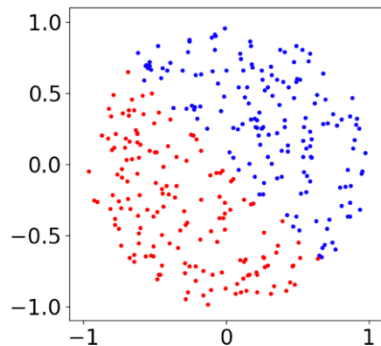
- Suppose
 - D is the training set, and its compromised version is \tilde{D}
 - Differentially-private learner: M
- Goals
 - Minimize the objective function: $J(\tilde{D}) := \mathbf{E}_b \left[C(\mathcal{M}(\tilde{D}, b)) \right]$
 - **Three attacks**
 - Parameter-targeting attack: make the model $\tilde{\theta}$ to *be close to* a target θ
 - Label-targeting attack: cause *small* prediction error on $\{z_i^*\}_{i \in [m]}$
 - Label-aversion attack: induce *large* prediction error on $\{z_i^*\}_{i \in [m]}$
- Capability
 - Modify k items in D

TRAINING-TIME DEFENSES: DIFFERENTIAL PRIVACY

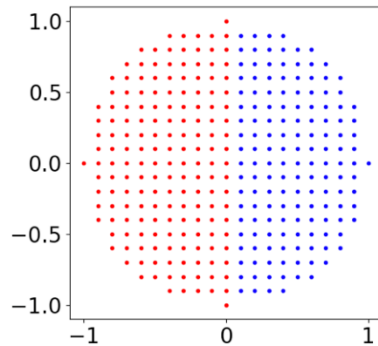
- DP as a poisoning defense
 - Construct the lower-bound $J(\tilde{D}) \geq e^{-k\epsilon} J(D)$
- One-shot kill attack (single-poison attack)
 - $k = 1$: the lower bound becomes $J(\tilde{D}) \geq e^{-\epsilon} J(D)$
 - $k \geq \lceil 1/\epsilon \log \tau \rceil$ modification can achieve $J(\tilde{D}) \geq 1/\tau J(D)$
 - ...

EVALUATION

- Setup [binary classification tasks]
 - Dataset: Synthetic data | Real data (UCI ML Repo.)
 - Models: Logistic regression | Ridge-regression
- Crafting poisons
 - Demonstrate on 2-D synthetic data



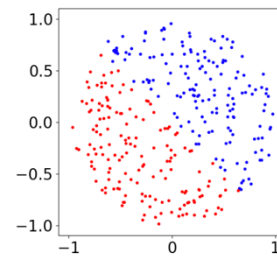
(b) training set



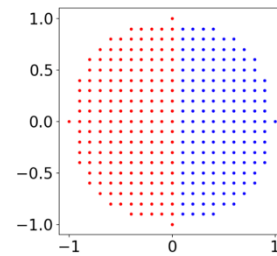
(c) evaluation set

EVALUATION

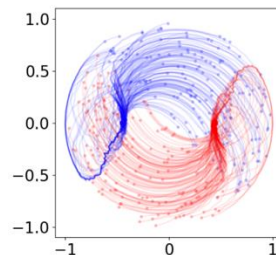
- Results of the three attacks on 2-D artificial data
 - Set $k = n$
 - Each attack achieves its objective



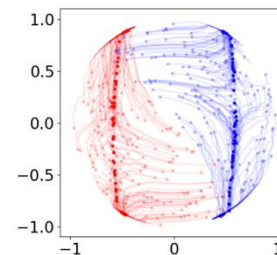
(b) training set



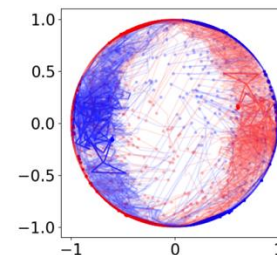
(c) evaluation set



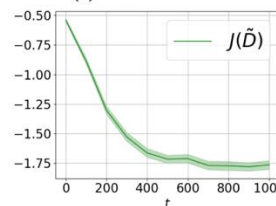
(a) label-aversion



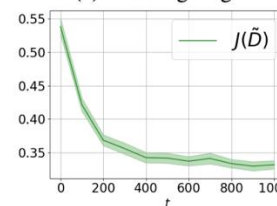
(b) label-targeting



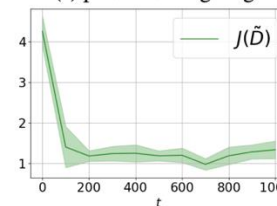
(c) parameter-targeting



(d) label-aversion



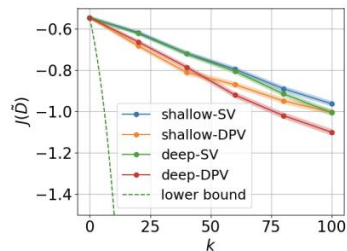
(e) label-targeting



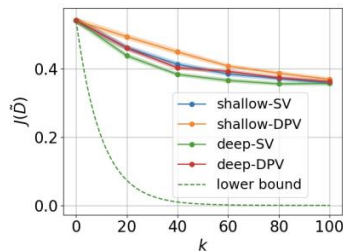
(f) parameter-targeting

EVALUATION

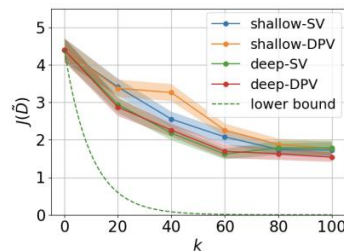
- Results of the three attacks on 2-D artificial data
 - The attack cost decreases as k increases (the attack becomes easier!)



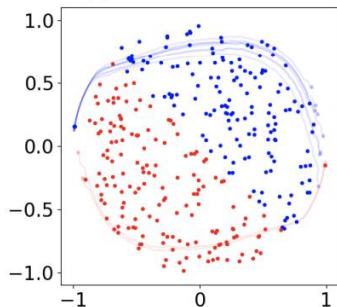
(a) label-aversion



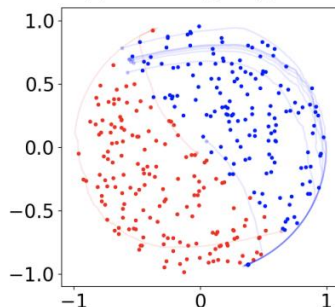
(b) label-targeting



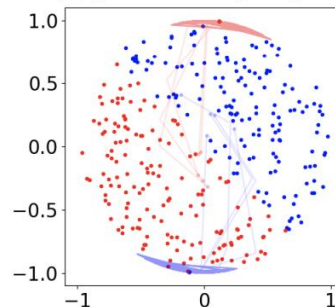
(c) parameter-targeting



(d) label-aversion



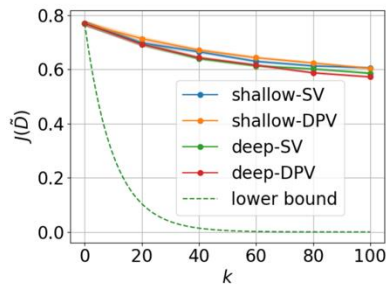
(e) label-targeting



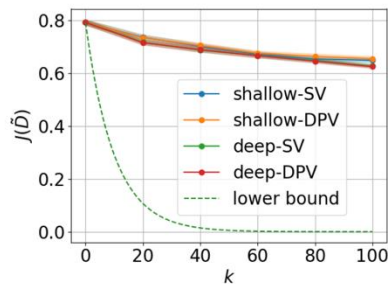
(f) parameter-targeting

EVALUATION

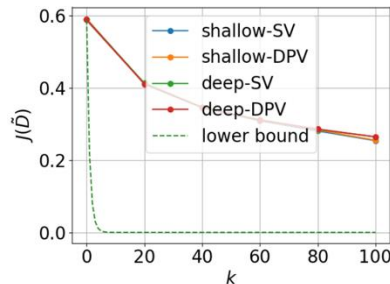
- Results of the *label-targeting* attacks on real-world datasets
 - (left) vs. logistic regression, (right) vs. ridge regression
 - The attacks work well also on the DP learners
 - The gap between the lower bound and the actual attack success exists



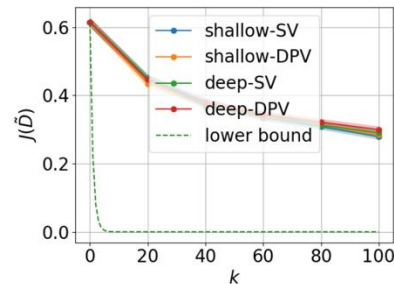
(a) objective perturbation



(b) output perturbation



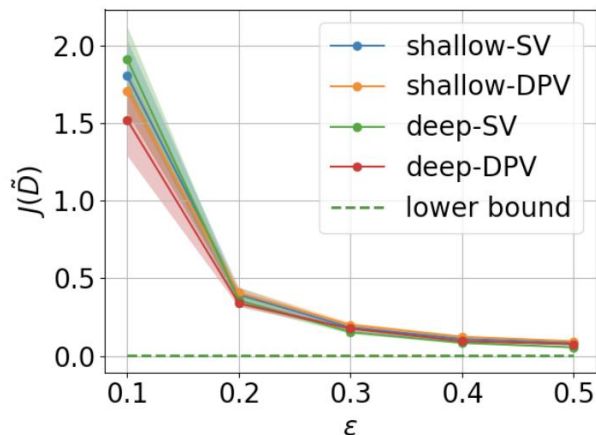
(a) objective perturbation



(b) output perturbation

EVALUATION

- Results of the *label-targeting* attacks on real-world datasets
 - In DP, the attack costs significantly higher than the case w/o DP
 - ex. with 20 poisons, the cost **w/o DP** is almost **zero** whereas **with DP, it's 0.4**
- Interesting Observation!
 - Attacks are much easier with weak (small epsilon) privacy



Thank You!

Sanghyun Hong

<https://secure-ai.systems/courses/MLSec/current>



Oregon State
University

SAIL
Secure AI Systems Lab