Robust Learning for Data Poisoning Attacks

Background

- Machine Learning systems are fragile, susceptible to attacks.
- Types of attacks: inference-time attacks, data poisoning attacks.
- Data poisoning attacks: the adversary manipulates the training data.
- $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ $\tilde{S} = \{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^n$ $\tilde{S} = \{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^n$
- Backdoor attack,
- Clean label attack: $\tilde{\mathbf{x}}_i = \mathbf{x}_i + \delta_i$, $\tilde{y}_i = y_i$,
- Label flip attack: $\tilde{x}_i = x_i$, $\tilde{y}_i = -y_i$ with probability β .

Notation: $B = \max \|\delta_i\|_2$ is per-sample perturbation; S =

is overall perturbation; β is probability of label flips; $\tilde{\mathcal{O}}$ hides polylogarithmic dependence on n.

Convex Learning Problem (Warm-up)

- Input/Label spaces $\mathscr{X} \subseteq \mathbb{R}^d$, $\mathscr{Y} = \{\pm 1\}$; distribution \mathscr{D}
- Goal: solve the stochastic optimization problem

$$\min_{\mathbf{w}\in \mathsf{W}} F(\mathbf{w}) := \mathbb{E}_{(\mathbf{x}, \mathbf{y})\sim \mathscr{D}}[\ell(\mathbf{y}f(\mathbf{x}; \mathbf{w}))],$$

where W is a convex set, ℓ is convex in W.

 Standard approach is to use SGD, where the learner tak gets access to a first order stochastic oracle for $\hat{g}(w) \in$ Observation: data poisoning attacks (δ_i) can be viewed as poisoning attacks (ζ_i).

•
$$\delta_i = \tilde{\mathbf{x}}_i - \mathbf{x}_i$$
.

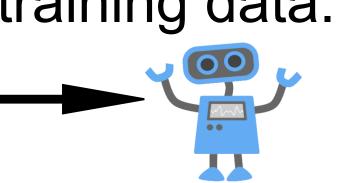
•
$$\zeta_i = \tilde{g}(w_i) - \hat{g}(w_i).$$

Theorem (Robustness of SGD). Excess risk bound for clean label attacks:

$$\mathbb{E}[F(\bar{\mathbf{w}})] - F(\mathbf{w}_*) \le O(\frac{1}{\sqrt{n}} + \frac{\sum_{i < n} \|\zeta_i\|}{n})$$

where the expectation is w.r.t. the initialization and the train Remark: 1. $\sum \|\zeta_i\| = \mathcal{O}(\sqrt{n})$ gives no significant statistic l < n

2. The above upper bound is tight in an information sense (see paper for a lower bound).



$$= \sum_{i=1}^{n} ||\delta_i||_2$$

es poly-

$$\mathcal{D} \text{ on } \mathcal{X} \times \mathcal{Y}.$$

kes w and

$$\in \partial F(w)$$
.

Two-layer neural networks

A two-layer ReLU net parameterized by (a, W

$$f(x; a, W) := \frac{1}{\sqrt{m}} \sum_{s=1}^{m} a_s \sigma(w_s^T x)$$
, ReLU: $\sigma(z)$, ne

Trained by online SGD using logistic loss.

Goal: minimize $L(W) := \mathbb{P}_{(x,y)\sim \mathcal{D}}(yf(x;a,W) < 0).$

Key Assumptions

There exists a margin parameter $\gamma > 0$, and a linear separator $\bar{v}: \mathbb{R}^d \to \mathbb{R}^d$, s.t.

- $\mathbb{E}_{z}[\|\bar{v}(z)\|^{2}] < \infty$,
- $\|\bar{v}(z)\|_2 \leq 1$ for all $z \in \mathbb{R}^d$,
- $\mathbb{E}_{z \sim \mathcal{N}(0, I_d)}[y\langle \overline{v}(z), x \mathbb{I}[z^\top x \ge 0] \rangle] \ge \gamma$ for almost all $(x, y) \sim \mathcal{D}$.

Main Results

Theorem: With probability at least $1 - \delta$, we show the following for the iterates of SGD, Regime A (clean label attack, large per-sample perturbation, small overall perturbation): $\frac{1}{2}\sum L(W_i) \lesssim \frac{\ln^2(\sqrt{n/4}) + \ln(24n/\delta)}{2}$

provided that $B \leq \tilde{\mathcal{O}}(\gamma/\sqrt{d}), \ \tilde{\mathcal{O}}(1/\gamma^8) \leq m \leq \tilde{\mathcal{O}}(\gamma/\sqrt{d})$ Remark: Regime A requires $S \lesssim \gamma^2 \sqrt{n}$ to allow a non-empty width range. Regime B (clean label attack, small per-sample perturbation, large overall perturbation): $\ln^2(\sqrt{n/4}) + \ln(24n/\delta)$

$$\frac{1}{n} \sum_{i < n} L(\mathbf{W}_i) \lesssim \cdot$$

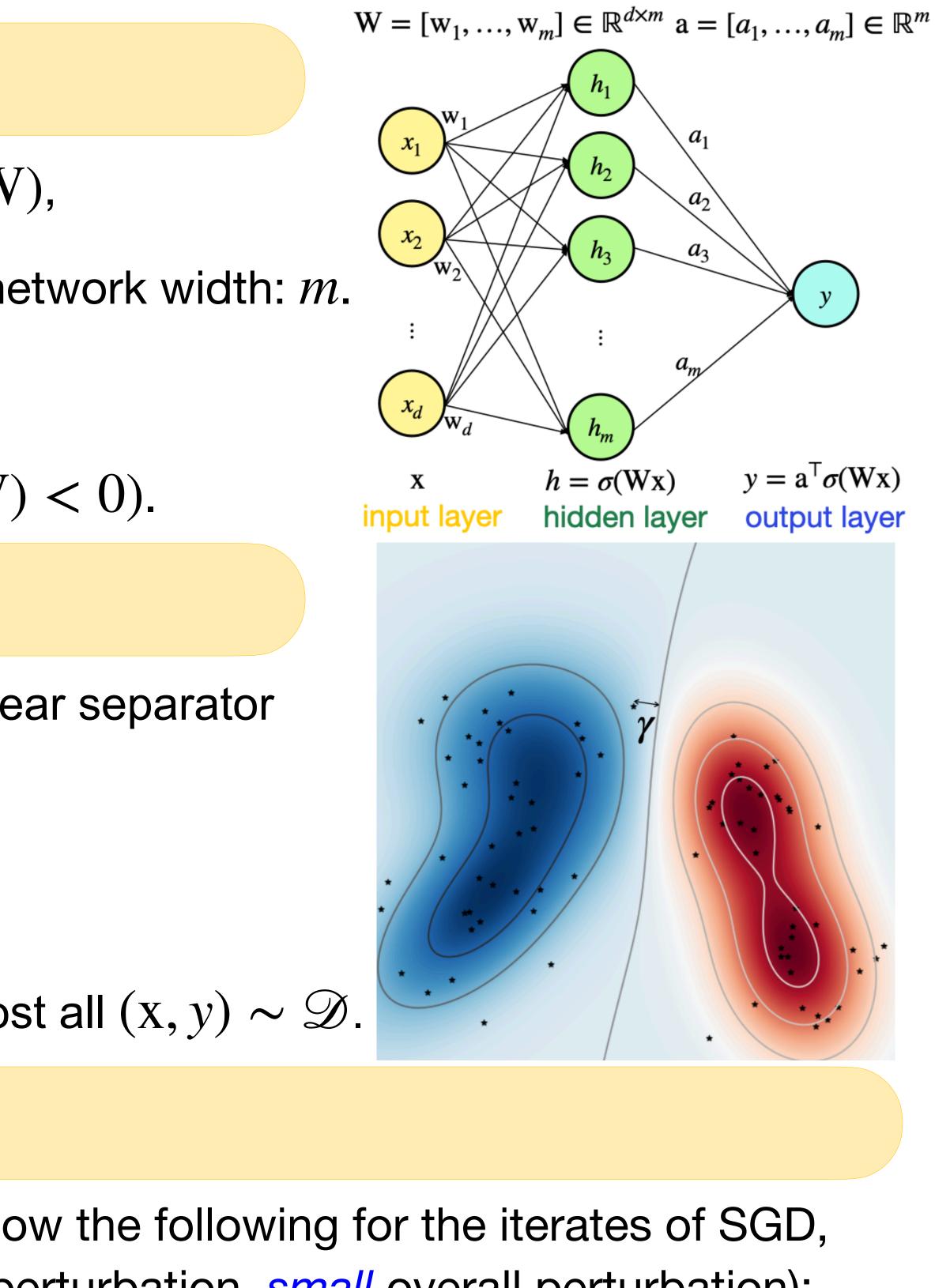
for $m \ge \tilde{\mathcal{O}}(1/\gamma^8)$, provided that $B \lesssim \min\{\frac{1}{\sqrt{md}}$

Remark: Regime B allows a larger overall perturb Regime C (label flip attack):

ning complee	$\frac{1}{n} \sum_{i < n} L(\mathbf{W}_i) \lesssim \frac{\ln n}{n}$
aining samples.	$2 (1 \cdot 8)$ $\ln(n/$
stical overhead.	for $m \geq \tilde{\mathcal{O}}(1/\gamma^8)$, provided that $\beta \lesssim \frac{\ln(n/\delta)}{(\sqrt{\ln(n/\delta)} + \delta)^2}$
tion-theoretic	Remark: 1. For regime A and C, the generalization
	the effective overall perturbation budg
	2. All three regimes require an upperbour

Yunjuan Wang, Poorya Mianjy, Raman Arora

Department of Computer Science, Johns Hopkins University



$$\sqrt{n\gamma^2}$$

 $\delta(n/\gamma^4 S^2).$

$$n\gamma^2$$
 γ γ

$$+\sqrt{m\ln(m/\delta)}, \gamma + \sqrt{d} + \sqrt{\ln(mn/\delta)}$$

pation budget of order $\mathcal{O}(n)$.

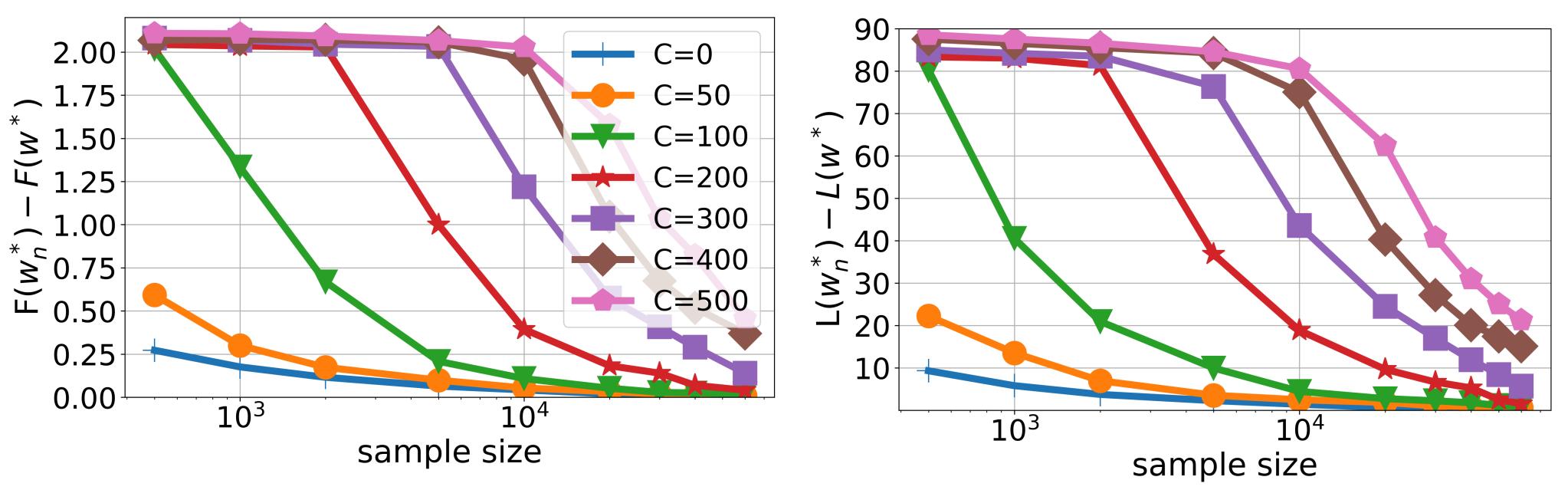
$$\frac{\ln^{2}(\sqrt{n/4}) + \ln(16n/\delta)}{\sqrt{n\gamma^{2}}}$$

$$\frac{n/\delta}{\sqrt{n}} + \ln^{2}(n)}{\sqrt{m(1+\gamma)}\gamma^{2}\sqrt{n}}.$$

on bounds are of the same rate of $\mathcal{O}(1/\sqrt{n})$, dget are almost of the same order $\tilde{O}(\sqrt{n})$. 2. All three regimes require an upperbound and lowerbound on the network width.

Experiments

- Generate the poisoned data:



on the given sample of size *n*.

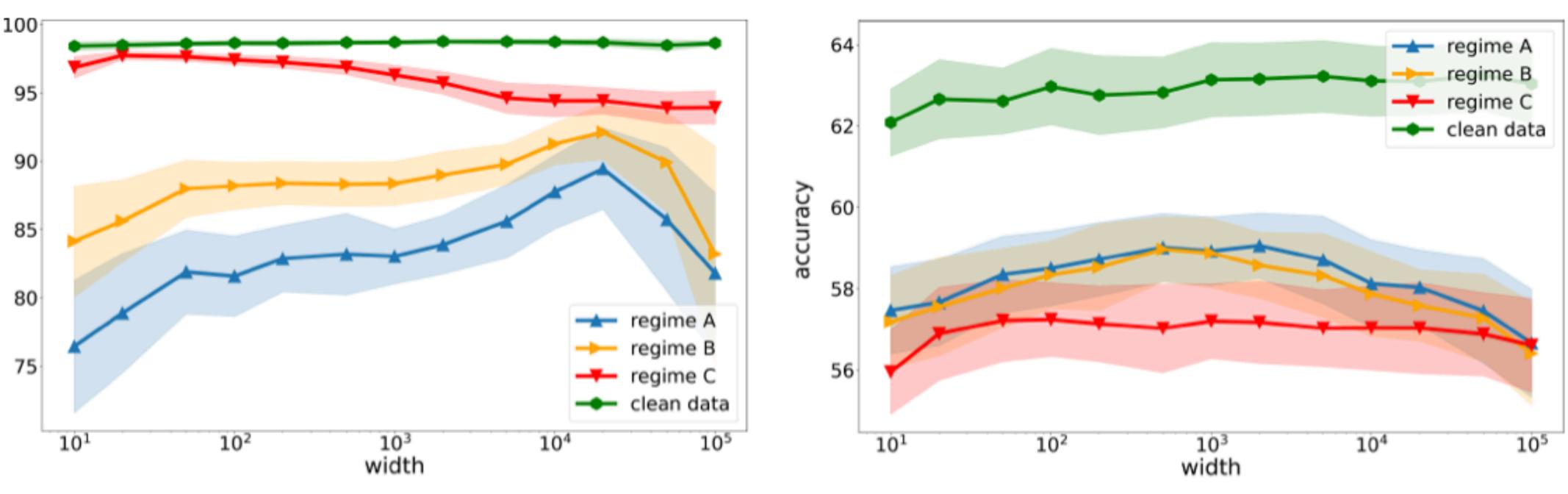


Figure: Clean test accuracy as a function of network width under clean data setting and poisoned data setting on MNIST (left) and CIFAR10 (right).

- they're too wide.

Main Takeaway: networks that are extremely over-parameterized are more susceptible to attacks.



 Propose negated loss by flipping the sign on both the model prediction and the cross-entropy loss, i.e $\ell'_{-}(z) := -\ell'(-z)$, where z is the model prediction.

1) Use mini-batch SGD to learn the best model parameters W^* on the clean data. 2) Take a stochastic gradient ascent step on the negated loss to maximize the negated loss function $\ell'_{(x;w^*)}$ with respect to x.

3) Project onto the constraints -- the $\ell_{2,1}$ (regime A) or $\ell_{2,\infty}$ -norm ball (regime B).

Figure: The excess loss $F(w_n^*) - F(w^*)$ (left); and the excess error $L(w_n^*) - L(w^*)$ (right), as a function of sample size n with different corruption parameter C under regime A. Specify overall budget $S = C\sqrt{n}$. Here, w_n^* denotes the optimal parameters

• The generalization accuracy decreases if the models are not wide enough or if

• The inverted U curve challenges the nascent view in the deep learning literature that larger models generalize better, at least under adversarial perturbation.