
Certifying robustness to adaptive data poisoning

Avinandan Bose
University of Washington
avibose@cs.washington.edu

Madeleine Udell
Stanford University
udell@stanford.edu

Laurent Lessard
Northeastern University
l.lessard@northeastern.edu

Maryam Fazel
University of Washington
mfazel@uw.edu

Krishnamurthy Dj Dvijotham
ServiceNow Research
dvij@cs.washington.edu

Abstract

The rise of foundation models fine-tuned on human feedback from potentially untrusted users has increased the risk of adversarial data poisoning, necessitating the study of robustness of learning algorithms against such attacks. Existing research on provable certified robustness against data poisoning attacks primarily focuses on certifying robustness for static adversaries that modify a fraction of the dataset used to train the model *before the training algorithm is applied*. In practice, particularly when learning from human feedback in an online sense, adversaries can observe and react to the learning process and inject poisoned samples that optimize adversarial objectives better than when they are restricted to poisoning a static dataset once before the learning algorithm is applied. Indeed, it has been shown in prior work that online dynamic adversaries can be significantly more powerful than static ones. We present a novel framework for computing certified bounds on the impact of dynamic poisoning, and use these certificates to design robust learning algorithms. We give an illustration of the framework for the mean-estimation problem and outline directions for extending this in further work.

1 Introduction & Problem Formulation

With the advent of foundation models fine tuned using human feedback gathered from potentially untrusted users (for example, users of a publicly available language model [4, 11]), the potential for adversarial or malicious data entering the training data of a model increases substantially. This motivates the study of robustness of learning algorithms to poisoning attacks [1]. More recently, there have been works that attempt to achieve “certified robustness” to data poisoning, i.e., proving that the worst case impact of poisoning is below a certain bound that depends on parameters of the learning algorithm. All the work in this space, to the best of our knowledge, focuses on the *static* poisoning adversary [15, 21]. Even in [17] which is the closest setting to our work, the poisoning adversary acts over offline datasets in a temporally extended fashion which are poisoned in one shot, and thus is not dynamic. There has been work on *dynamic* attack algorithms [20, 18] showing that these attacks can indeed be more powerful than static attacks. This motivates the question we study: can we obtain certificates of robustness for a broad class of learning algorithms against *dynamic* poisoning adversaries?

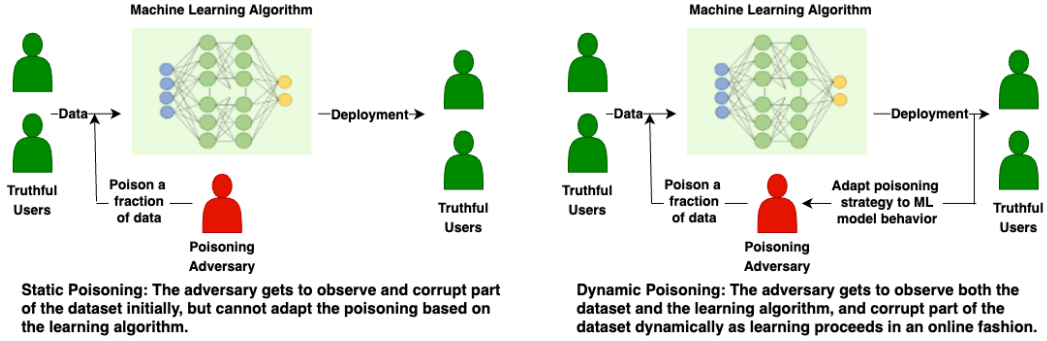


Figure 1: A schematic diagram to highlight the differences between static and dynamic poisoning.

Attack Type	Adversary adapts poisoning strategy upon observing model behavior	Adversary can poison data for deployed model	Certified robustness
Static / One-shot ([16, 15, 13])	✗	✗	✓
Backdoor([3, 7, 8, 22, 21, 14])	✗	✓	✓
Dynamic attack only ([18, 20])	✓	✓	✗
Dynamic attack & defense (Ours)	✓	✓	✓

Table 1: A comparison with lines of work closest to ours. A detailed description is provided in Appendix A.

In this paper, we study learning algorithms corrupted by a dynamic poisoning adversary who can observe the behavior of the learning algorithm and adapt the poisoning in response. This is relevant in scenarios where models are continuously/periodically updated in the face of new feedback, as is common in RLHF/fine tuning applications. We provide (to the best of our knowledge) the first general framework for computing certified bounds on the worst case impact of a data poisoning attacker, and further, use this certificate to design robust learning algorithms. We give an illustration of the framework for the mean-estimation problem (see Section 3) and suggest directions of future work to apply the framework to more realistic learning scenarios.

2 Problem setup

We now develop the exact problem setup that we study in the paper. We will assume that the learning algorithm we study is aimed at estimating parameters $\theta \in \Theta$ and each step of the learning algorithm makes updates to the estimate of these parameters based on potentially poisoned data. The following components fully define the problem setup.

Online learning algorithm We consider learning algorithms that operate online receiving at each step a new datapoint and making an update to parameters being estimated. In

Notation	Interpretation	Belongs to
θ	Parameters of model	Θ
ϕ	Hyper-parameters of learning algorithm	Φ
w	Gaussian noise injected into learning algorithm	\mathcal{W}
z	Datapoint	\mathcal{Z}
F	Update rule of learning algorithm	$\Theta \times \mathcal{W} \times \mathcal{Z} \mapsto \Theta$
z^{adv}	Adversarial data point	\mathcal{A}

Table 2: Notation

particular, we consider learning algorithms that can be written as

$$\theta_{t+1} \leftarrow F \left(\underbrace{\theta_t}_{\text{Parameter estimate at time } t}, \underbrace{w_t}_{\text{Exogeneous noise input}}, \underbrace{z_t}_{\text{Datapoint received at time } t} \right) \quad (1)$$

where $F : \Theta \times \mathcal{W} \times \mathcal{Z} \mapsto \Theta$ is an update function that maps the parameters at time t to new parameters, given an exogeneous noise input w_t and a datapoint z_t . The exogenous noise input refers to noise artificially injected into the training algorithm in order to make the algorithm more robust to potential poisoning. We further assume that the distribution of w_t is independent of t and each w_t is sampled iid.

Poisoned learning algorithm We work in a setting where some of the datapoints received by the learning algorithm are corrupted by an adversary, with the corruption allowed to be a function of the entire trajectory of the learning algorithm up to that point. We assume that with a fixed probability the data point the algorithm receives at each time step is poisoned. In practice, this could reflect the situation that out of a large population of human users providing feedback to a learning system, a small fraction are adversarial and will provide poisoned feedback.

Mathematically, we have that at time t the learning algorithm receives a datapoint $z_t \sim \epsilon \text{Dirac}(z_t^{\text{adv}}) + (1 - \epsilon) \mathbb{P}^{\text{data}}$ and ϵ is a parameter that controls the “level” of poisoning (analogous to the fraction of poisoned samples in static poisoning settings [15]). This is a special case of Huber’s contamination model, which is used in the robust statistics literature [5] (with the contamination model being a Dirac distribution). We restrict the adversary to choose $z_t^{\text{adv}} \in \mathcal{A}$ which reflects the allowed range of datapoints due to input feature normalization or outlier detection systems.

Adversarial objective We assume that the poisoning adversary is interested in maximizing some adversarial objective $\ell_{\text{adv}} : \Theta \mapsto \mathbb{R}$, for example, the expected prediction error on some target distribution of interest to the adversary.

Dynamics as a Markov Chain The dynamics (1) gives rise to a Markov chain over the parameters θ . If \mathbb{P}_t denotes the distribution over parameters at time t , we have

$$\mathbb{P}_{t+1}(\theta) = \int \mathbb{P}_{F, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\theta | \theta') \mathbb{P}_t(\theta') d\theta',$$

where $\mathbb{P}_{F, \mathbb{P}^{\text{data}}, z^{\text{adv}}}$ is the transition kernel induced, i.e., the conditional probability distribution of $\theta' = F(\theta, w_t, z_t)$ by (1) given θ .

2.1 Technical Approach: Certificate of Robustness

We are now ready to present our technical result, a certificate of robustness against dynamic data poisoning adversaries. Since the learning algorithm is a Markov process, the optimal sequence of actions for the adversary (i.e., choices of z^{adv}) constitute a Markov Decision Process with

$$\text{States } \theta, \text{ Actions } z^{\text{adv}}, \text{ Transition Kernel } \mathbb{P}^{\text{trans}}(\theta' | \theta, z^{\text{adv}}) = \mathbb{P}_{F, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\theta' | \theta),$$

and hence, can be formulated as an *infinite dimensional* linear program [12]. In particular, for the infinite horizon average reward setting [9], the LP can be written as

$$\sup_{\mathbb{P} \in \mathcal{P}[\Theta \times \mathcal{Z}]} \mathbb{E}_{\theta, z^{\text{adv}} \sim \mathbb{P}} [\ell_{\text{adv}}(\theta)] \quad (2a)$$

$$\text{subject to } \mathbb{E}_{\theta, z^{\text{adv}} \sim \mathbb{P}} \left[\mathbb{P}_{F, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\theta' | \theta) \right] = \mathbb{E}_{\theta, z^{\text{adv}} \sim \mathbb{P}} [\mathbb{I}[\theta' = \theta]] \quad \forall \theta' \in \mathbb{R}^d. \quad (2b)$$

where $\mathcal{P}[\Theta \times \mathcal{Z}]$ denotes the space of probability measures on $\Theta \times \mathcal{Z}$ and \mathbb{I} denotes the indicator function that equals 1 if its argument is true and 0 otherwise.

Theorem 1. For any function $\lambda : \Theta \mapsto \mathbb{R}$, we can upper bound the optimal value of (2) by

$$\sup_{\substack{\theta \in \Theta \\ z^{\text{adv}} \in \mathcal{A}}} \mathbb{E}_{\theta' \sim \mathbb{P}_{F, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\cdot | \theta)} [\lambda(\theta')] + \ell_{\text{adv}}(\theta) - \lambda(\theta). \quad (3)$$

Proof. Follows by weak duality for the LP (4a).

If strong duality holds, we further have that the optimal value of (2) is exactly equal to

$$\inf_{\lambda: \Theta \mapsto \mathbb{R}} \sup_{\theta \in \Theta, z^{\text{adv}} \in \mathcal{A}} \mathbb{E}_{\theta' \sim \mathbb{P}_{F, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\cdot | \theta)} [\lambda(\theta')] + \ell_{\text{adv}}(\theta) - \lambda(\theta). \quad (4a)$$

2.2 Meta-learning a robust learning algorithm

Consider a parameterized family of learning algorithms F_ϕ with tunable parameters $\phi \in \Phi$. Based on the above analysis, we can attempt to design the parameters ϕ of the learning algorithm to trade-off performance and robustness. In particular, in the absence of poisoned data, assume that the updates (1) result in a stationary distribution $\mathbb{P}(\phi, \mathbb{P}^{\text{data}})$ over model parameters θ .

Given some space of data distributions \mathcal{P} we can sample from (in a meta learning sense), we can propose the following criterion:

$$\inf_{\substack{\phi \in \Phi \\ \lambda: \Theta \mapsto \mathbb{R}}} \mathbb{E}_{\mathbb{P}^{\text{data}} \sim \mathcal{P}} \left[\mathbb{E}_{\theta \sim \mathbb{P}(\phi, \mathbb{P}^{\text{data}})} [\ell(\theta)] + \kappa \left(\sup_{\theta \in \Theta, z^{\text{adv}} \in \mathcal{A}} \mathbb{E}_{\theta' \sim \mathbb{P}_{S, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\cdot | \theta)} [\lambda(\theta')] + \ell_{\text{adv}}(\theta) - \lambda(\theta) \right) \right], \quad (5)$$

where $\kappa > 0$ is a trade-off parameter. The outer expectation is a meta-learning inspired formulation, where we are designing a learning algorithm that is good "in expectation" under a meta-distribution over distributions. The first term in the outer expectation constitutes "doing well" in the absence of the adversary by converging to a stationary distribution over parameters that incurs low expected loss. The second term is an upper bound on the worse case loss incurred by the learning algorithm in the presence of the adversary.

3 Mean estimation

Consider the mean estimation problem, where we aim to learn the parameter $\theta \in \mathbb{R}^d$ to estimate the mean $\mu = \mathbb{E}_{z \sim \mathbb{P}^{\text{data}}} [z]$ of a distribution \mathbb{P}^{data} . Given a data point z_t , the learning rule is given by:

$$\theta_{t+1} \leftarrow (1 - \eta)\theta_t + \eta z_t + \eta \mathbf{B} w_t,$$

where $S = \mathbf{B}\mathbf{B}^\top \in S_+^d$ is the tunable defense parameter and $w_t \sim \mathcal{N}(0, \mathbf{I})$ is Gaussian noise. The adversarial loss is given by:

$$\ell_{\text{adv}}(\theta) = \|\mu - \theta\|^2.$$

Certificate on adversarial loss (analysis)

Theorem 2. Choosing $\lambda : \mathbb{R}^d \rightarrow \mathbb{R}$ in Theorem 1 to be quadratic, i.e. $\lambda(\theta) = \theta^\top \mathbf{A} \theta + \theta^\top \mathbf{b}$, the adversarial constraint set of the form $\|z^{\text{adv}} - \mu\|_2^2 \leq r$, the certificate for the mean estimation problem for $\mathbb{P}^{\text{data}}(z) = \mathcal{N}(z | \mu, \Sigma)$ for a fixed learning algorithm (i.e. S is fixed) is given by:

$$\inf_{A \in S^d, \mathbf{b} \in \mathbb{R}^d, \nu \geq 0} g(\mathbf{A}, \mathbf{b}, \nu, S, \mu, \Sigma), \quad (6)$$

where $g(\mathbf{A}, \mathbf{b}, \nu, S, \mu, \Sigma)$ is a convex objective in $\mathbf{A}, \mathbf{b}, \nu$ (matrix fractional objective with Linear Matrix Inequality (LMI) constraint) as defined below:

$$g(\mathbf{A}, \mathbf{b}, \nu, S, \mu, \Sigma) = \begin{cases} \frac{1}{4} \left\| \begin{bmatrix} 2(1 - \epsilon)\eta(1 - \eta)\mathbf{A}\mu - 2\mu - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} + 2\nu\mu \end{bmatrix} \right\|_D^2 + (1 - \epsilon)(\eta^2 \text{Trace}(\Sigma\mathbf{A}) + \eta^2 \mu^\top \mathbf{A} \mu + \eta \mathbf{b}^\top \mu) \\ + \mu^\top \mu + \eta^2 \text{Trace}(\mathbf{A}S) + \nu(r - \mu^\top \mu) \text{ if } \nu \geq 0; \mathbf{D} \succeq 0 \\ -\infty \text{ else} \end{cases} \quad (7)$$

Algorithm 1 Meta learning a robust learning algorithm for mean estimation

- 1: **Input:** Set of K distributions $\{\mathcal{N}(\mu_i, \Sigma_i)\}_{i \in [K]}$ sampled from \mathcal{P} , tradeoff parameter κ .
 - 2: **Initialize:** $\mathbf{S} \in \mathbb{S}_+^d$ randomly.
 - 3: *Alternating Minimization over Lagrange multipliers $\{\mathbf{A}_i, \mathbf{b}_i, \nu_i\}_{i \in [K]}$ and defence parameter \mathbf{S} .*
 - 4: **for** $t = 1, \dots$ **do**
 - 5: **for** $i = 1, \dots, K$ **do**
 - 6: $\mathbf{A}_i, \mathbf{b}_i, \nu_i = \inf_{\mathbf{A} \in \mathbb{S}^d, \mathbf{b} \in \mathbb{R}^d, \nu \geq 0} g(\mathbf{A}, \mathbf{b}, \nu, \mathbf{S}, \mu_i, \Sigma_i)$.
 - 7: **end for**
 - 8: $\mathbf{S} = \operatorname{argmin}_{\mathbf{S} \in \mathbb{S}_+^d} (\eta^2 \operatorname{Trace}(\mathbf{S}) + \frac{\kappa}{K} \sum_{i \in [K]} g(\mathbf{A}_i, \mathbf{b}_i, \nu_i, \mathbf{S}, \mu_i, \Sigma_i))$.
 - 9: **end for**
-

and $\mathbf{D} = \begin{bmatrix} (1 - (1 - \eta)^2)\mathbf{A} - \mathbf{I} & -\eta\epsilon(1 - \eta)\mathbf{A} \\ -\eta\epsilon(1 - \eta)\mathbf{A} & -\epsilon\eta^2\mathbf{A} - \nu\mathbf{I} \end{bmatrix}$ and $\|\mathbf{x}\|_{\mathbf{D}}^2 = \mathbf{x}^\top \mathbf{D}^{-1} \mathbf{x}$.

Meta-Learning Algorithm Following the formulation in Eq. (5), we wish to learn a defense parameter \mathbf{S} that minimizes the expected loss (expectation over different \mathbb{P}^{data} from the meta distribution \mathcal{P}). For the mean estimation problem this boils down to solving:

$$\inf_{\mathbf{S} \in \mathbb{S}_+^d} \eta^2 \operatorname{Trace}(\mathbf{S}) + \kappa \mathbb{E}_{\mu, \Sigma \sim \mathcal{P}} \left[\inf_{\substack{\nu \geq 0 \\ \mathbf{A} \in \mathbb{S}^d, \mathbf{b} \in \mathbb{R}^d}} g(\mathbf{A}, \mathbf{b}, \nu, \mathbf{S}, \mu, \Sigma) \right]. \quad (8)$$

In practice, one observes a finite number of distributions from \mathcal{P} , and sample average approximation is leveraged, with the aim of learning a defense parameter which generalizes well to unseen distributions from \mathcal{P} . This process is stated in Algorithm 1.

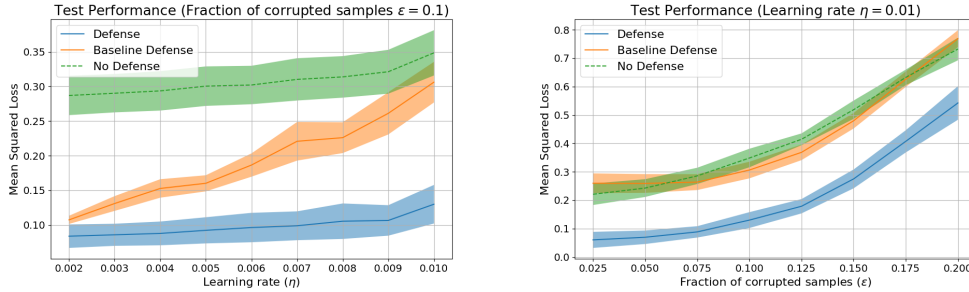


Figure 2: Test performance (mean squared error between true and estimated means) on 50 ($d = 20$ dimensional) Gaussian distributions drawn from Gaussian prior for the mean and Inverse Wishart prior for the covariance. The defense parameter \mathbf{S} was trained with 10 such randomly chosen Gaussians via Algorithm 1. We varied the learning rates (left) and the the fraction of samples corrupted by the dynamic adversary (right) and observe that our defense beats training without defense $\epsilon = 0.1$ significantly.

References

- [1] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. *arXiv preprint arXiv:1206.6389*, 2012.
- [2] Battista Biggio and Fabio Roli. Wild patterns: Ten years after the rise of adversarial machine learning. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pages 2154–2156, 2018.
- [3] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- [4] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- [5] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge university press, 2023.
- [6] Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- [7] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- [8] Xingshuo Han, Guowen Xu, Yuan Zhou, Xuehuan Yang, Jiwei Li, and Tianwei Zhang. Physical backdoor attacks to lane detection systems in autonomous driving. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 2957–2968, 2022.
- [9] Alan Malek, Yasin Abbasi-Yadkori, and Peter Bartlett. Linear programming for large-scale markov decision problems. In *International conference on machine learning*, pages 496–504. PMLR, 2014.
- [10] James Newsome, Brad Karp, and Dawn Song. Paragraph: Thwarting signature learning by training maliciously. In *Recent Advances in Intrusion Detection: 9th International Symposium, RAID 2006 Hamburg, Germany, September 20-22, 2006 Proceedings 9*, pages 81–105. Springer, 2006.
- [11] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [12] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [13] Elan Rosenfeld, Ezra Winston, Pradeep Ravikumar, and Zico Kolter. Certified robustness to label-flipping attacks via randomized smoothing. In *International Conference on Machine Learning*, pages 8230–8241. PMLR, 2020.
- [14] Philip Sosnin, Mark N Müller, Maximilian Baader, Calvin Tsay, and Matthew Wicker. Certified robustness to data poisoning in gradient-based training. *arXiv preprint arXiv:2406.05670*, 2024.
- [15] Jacob Steinhardt, Pang Wei W Koh, and Percy S Liang. Certified defenses for data poisoning attacks. *Advances in neural information processing systems*, 30, 2017.
- [16] Zhiyi Tian, Lei Cui, Jie Liang, and Shui Yu. A comprehensive survey on poisoning attacks and countermeasures in machine learning. *ACM Computing Surveys*, 55(8):1–35, 2022.
- [17] Wenxiao Wang and Soheil Feizi. Temporal robustness against data poisoning. *Advances in Neural Information Processing Systems*, 36, 2024.

- [18] Yizhen Wang and Kamalika Chaudhuri. Data poisoning attacks against online learning. *arXiv preprint arXiv:1808.08994*, 2018.
- [19] Chulin Xie, Yunhui Long, Pin-Yu Chen, and Bo Li. Uncovering the connection between differential privacy and certified robustness of federated learning against poisoning attacks. *arXiv preprint arXiv:2209.04030*, 2022.
- [20] Xuezhou Zhang, Xiaojin Zhu, and Laurent Lessard. Online data poisoning attacks. In *Learning for Dynamics and Control*, pages 201–210. PMLR, 2020.
- [21] Yuhao Zhang, Aws Albarghouthi, and Loris D’Antoni. Bagflip: A certified defense against data poisoning. *Advances in Neural Information Processing Systems*, 35:31474–31483, 2022.
- [22] Chen Zhu, W Ronny Huang, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein. Transferable clean-label poisoning attacks on deep neural nets. In *International conference on machine learning*, pages 7614–7623. PMLR, 2019.

A Related Work

Data Poisoning Modern machine learning pipelines involve training on massive, uncurated datasets that are potentially untrustworthy and of such scale that conducting rigorous quality checks becomes impractical. Poisoning attacks [1, 10, 2] pose big security concerns upon deployment of ML models. Depending on which stage (training / deployment) the poisoning takes place, they can be characterised as follows : 1. Static attacks : The model is trained on an offline dataset with poisoned data. Attacks could be untargeted, which aim to prevent training convergence rendering an unusable model and thus denial of service [16] or targeted which are more task-specific and instead of simply increasing loss, attacks of this kind seeks to make the model output wrong predictions on specific tasks. 2. Backdoor attacks: In this setting, the test / deployment time data can be altered [3, 7, 8, 22]. Attackers manipulate a small proportion of the data such that, when a specific pattern / trigger is seen at test-time, the model returns a specific, erroneous prediction. 3. Dynamic (and adaptive) attacks: In scenarios where models are continuously / periodically updated in the face of new feedback, as is common in RLHF / fine tuning applications, a dynamic poisoning adversary [18, 20] can observe the behavior of the learning algorithm and adapt the poisoning in response.

Certified Poisoning Defense Recently, there have been works that attempt to achieve ‘‘certified robustness’’ to data poisoning, i.e., proving that the worst case impact of *any* poisoning strategy is below a certain bound that depends on parameters of the learning algorithm. All the work in this space, to the best of our knowledge, focuses on the *static* or *backdoor* attack adversary. [15] provide certificates for linear models trained with gradient descent, [13] present a statistical upper-bound on the effectiveness of ℓ -2 perturbations on training labels for linear models using randomized smoothing, [21, 14] present a model-agnostic certified approach that can effectively defend against both trigger-less and backdoor attacks, [19] observe that differential privacy, which usually covers addition or removal of data points, can also provide statistical guarantees in some limited poisoning settings. Even in [17] which is the closest setting to our work, the poisoning adversary acts over offline datasets in a temporally extended fashion which are poisoned in one shot, and thus is not dynamic.

B Proofs

Theorem 2. *Choosing $\lambda : \mathbb{R}^d \rightarrow \mathbb{R}$ in Theorem 1 to be quadratic, i.e. $\lambda(\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{A} \boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{b}$, the adversarial constraint set of the form $\|\mathbf{z}^{adv} - \boldsymbol{\mu}\|_2^2 \leq r$, the certificate for the mean estimation problem for $\mathbb{P}^{\text{data}}(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for a fixed learning algorithm (i.e. \mathbf{S} is fixed) is given by:*

$$\inf_{\mathbf{A} \in \mathbb{S}^d, \mathbf{b} \in \mathbb{R}^d, \nu \geq 0} g(\mathbf{A}, \mathbf{b}, \nu, \mathbf{S}, \boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (6)$$

where $g(\mathbf{A}, \mathbf{b}, \nu, \mathbf{S}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a convex objective in $\mathbf{A}, \mathbf{b}, \nu$ (matrix fractional objective with Linear Matrix Inequality (LMI) constraint) as defined below:

$$g(\mathbf{A}, \mathbf{b}, \nu, \mathbf{S}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \begin{cases} \frac{1}{4} \left\| \begin{bmatrix} 2(1-\epsilon)\eta(1-\eta)\mathbf{A}\boldsymbol{\mu} - 2\boldsymbol{\mu} - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} + 2\nu\boldsymbol{\mu} \end{bmatrix} \right\|_{\mathbf{D}}^2 + (1-\epsilon)(\eta^2 \text{Trace}(\boldsymbol{\Sigma}\mathbf{A}) + \eta^2 \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} + \eta \mathbf{b}^\top \boldsymbol{\mu}) \\ + \boldsymbol{\mu}^\top \boldsymbol{\mu} + \eta^2 \text{Trace}(\mathbf{A}\mathbf{S}) + \nu(r - \boldsymbol{\mu}^\top \boldsymbol{\mu}) \text{ if } \nu \geq 0; \mathbf{D} \succeq 0 \\ -\infty \text{ else} \end{cases} \quad (7)$$

$$\text{and } \mathbf{D} = \begin{bmatrix} (1 - (1 - \eta)^2)\mathbf{A} - \mathbf{I} & -\eta\epsilon(1 - \eta)\mathbf{A} \\ -\eta\epsilon(1 - \eta)\mathbf{A} & -\epsilon\eta^2\mathbf{A} - \nu\mathbf{I} \end{bmatrix} \text{ and } \|\mathbf{x}\|_{\mathbf{D}}^2 = \mathbf{x}^\top \mathbf{D}^{-1} \mathbf{x}.$$

Proof. We can write the learning algorithm in Eq. (1) for the case of mean estimation as follows:

$$\boldsymbol{\theta}_{t+1} = F(\boldsymbol{\theta}_t, \mathbf{z}_t) + \eta \mathbf{B} \mathbf{w}_t,$$

where $F(\boldsymbol{\theta}, \mathbf{z}) = \boldsymbol{\theta}(1 - \eta) + \eta\mathbf{z}$, which is a linear transformation of $\boldsymbol{\theta}$ followed by additive Gaussian noise.

The transition distribution for the parameter is given by:

$$\mathbb{P}_{\mathbf{S}, \mathbb{P}^{\text{data}}, \mathbf{z}^{\text{adv}}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \epsilon \mathcal{N}(\boldsymbol{\theta}' | F(\boldsymbol{\theta}, \mathbf{z}^{\text{adv}}), \eta^2 \mathbf{S}) + (1 - \epsilon) \mathbb{E}_{\mathbf{z} \sim \mathbb{P}^{\text{data}}} \left[\mathcal{N}(\boldsymbol{\theta}' | F(\boldsymbol{\theta}, \mathbf{z}), \eta^2 \mathbf{S}) \right] \quad (9)$$

which is a Gaussian distribution whose mean depends linearly on $\boldsymbol{\theta}$ and \mathbf{z}^{adv} .

Then, we have from Eq. (4a) that the certified bound on the adversarial objective is given by:

$$\sup_{\boldsymbol{\theta}} \epsilon \mathbb{E}_{\mathbf{z}^{\text{adv}} \in \mathcal{A}} \mathbb{E}_{\boldsymbol{\theta}' \sim \mathcal{N}(F(\boldsymbol{\theta}, \mathbf{z}^{\text{adv}}), \eta^2 \mathbf{S})} [\lambda(\boldsymbol{\theta}')] + (1 - \epsilon) \mathbb{E}_{\mathbf{z} \sim \mathbb{P}^{\text{data}}} \left[\mathbb{E}_{\boldsymbol{\theta}' \sim \mathcal{N}(F(\boldsymbol{\theta}, \mathbf{z}), \eta^2 \mathbf{S})} [\lambda(\boldsymbol{\theta}')] \right] - \lambda(\boldsymbol{\theta}) + \ell_{\text{adv}}(\boldsymbol{\theta}) \quad (10a)$$

We choose $\lambda(\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{A} \boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{b}$ to be a quadratic function. Then we have:

$$= \sup_{\boldsymbol{\theta}} \epsilon \left(\lambda(F(\boldsymbol{\theta}, \mathbf{z}^{\text{adv}})) + \eta^2 \langle \nabla^2 \lambda(0), \mathbf{S} \rangle \right) + (1 - \epsilon) \mathbb{E}_{\mathbf{z} \sim \mathbb{P}^{\text{data}}} \left[\lambda(F(\boldsymbol{\theta}, \mathbf{z})) + \eta^2 \langle \nabla^2 \lambda(0), \mathbf{S} \rangle \right] - \lambda(\boldsymbol{\theta}) + \ell_{\text{adv}}(\boldsymbol{\theta}) \quad (10b)$$

$$= \sup_{\boldsymbol{\theta}} \left(- \left\| \begin{bmatrix} \boldsymbol{\theta} \\ \mathbf{z}^{\text{adv}} \end{bmatrix} \right\|_{E^{-1}}^2 + \begin{bmatrix} \boldsymbol{\theta} \\ \mathbf{z}^{\text{adv}} \end{bmatrix}^\top \begin{bmatrix} 2(1 - \epsilon)\eta(1 - \eta)\mathbf{A}\boldsymbol{\mu} - 2\boldsymbol{\mu} - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} \end{bmatrix} \right) + (1 - \epsilon)(\eta^2 \text{Trace}(\boldsymbol{\Sigma}\mathbf{A}) + \eta^2 \boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu} + \eta\mathbf{b}^\top \boldsymbol{\mu}) + \boldsymbol{\mu}^\top \boldsymbol{\mu}, \quad (10c)$$

$$\text{where } E = \begin{bmatrix} (1 - (1 - \eta)^2)\mathbf{A} - \mathbf{I} & -\eta\epsilon(1 - \eta)\mathbf{A} \\ -\eta\epsilon(1 - \eta)\mathbf{A} & -\epsilon\eta^2\mathbf{A} \end{bmatrix},$$

The dual function of this supremum (with dual variable ν) can be written as:

$$= \inf_{\nu \geq 0} \sup_{\boldsymbol{\theta}} \left(- \left\| \begin{bmatrix} \boldsymbol{\theta} \\ \mathbf{z}^{\text{adv}} \end{bmatrix} \right\|_{D^{-1}}^2 + \begin{bmatrix} \boldsymbol{\theta} \\ \mathbf{z}^{\text{adv}} \end{bmatrix}^\top \begin{bmatrix} 2(1 - \epsilon)\eta(1 - \eta)\mathbf{A}\boldsymbol{\mu} - 2\boldsymbol{\mu} - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} + 2\nu\boldsymbol{\mu} \end{bmatrix} \right) + (1 - \epsilon)(\eta^2 \text{Trace}(\boldsymbol{\Sigma}\mathbf{A}) + \eta^2 \boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu} + \eta\mathbf{b}^\top \boldsymbol{\mu}) + \boldsymbol{\mu}^\top \boldsymbol{\mu} + \nu(r - \boldsymbol{\mu}^\top \boldsymbol{\mu}) \quad (10d)$$

$$\text{where } D = \begin{bmatrix} (1 - (1 - \eta)^2)\mathbf{A} - \mathbf{I} & -\eta\epsilon(1 - \eta)\mathbf{A} \\ -\eta\epsilon(1 - \eta)\mathbf{A} & -\epsilon\eta^2\mathbf{A} - \nu\mathbf{I} \end{bmatrix}.$$

The inner supremum is a quadratic expression in $\mathbf{z}^{\text{adv}}, \boldsymbol{\theta}$. A finite supremum exists if the Hessian of the expression is negative semidefinite. Plugging in the tractable maximizer of the quadratic, we get:

$$\inf_{\nu \geq 0} \frac{1}{4} \left\| \begin{bmatrix} 2(1 - \epsilon)\eta(1 - \eta)\mathbf{A}\boldsymbol{\mu} - 2\boldsymbol{\mu} - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} + 2\nu\boldsymbol{\mu} \end{bmatrix} \right\|_D^2 + (1 - \epsilon)(\eta^2 \text{Trace}(\boldsymbol{\Sigma}\mathbf{A}) + \eta^2 \boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu} + \eta\mathbf{b}^\top \boldsymbol{\mu}) + \boldsymbol{\mu}^\top \boldsymbol{\mu} + \eta^2 \text{Trace}(\mathbf{A}\mathbf{S}) + \nu(r - \boldsymbol{\mu}^\top \boldsymbol{\mu}) \text{ such that } D \succeq 0. \quad (10e)$$

This completes the proof. □

Lemma B.1. *The stationary distribution in the absence of adversary in Eq. (??) for the mean estimation problem for $\mathbb{P}^{\text{data}} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ takes the form:*

$$\mathbb{P}(\mathbf{S}, \mathbb{P}^{\text{data}}) = \mathcal{N}(\boldsymbol{\mu}, \eta^2 \mathbf{S}).$$

Proof. The stationary distribution is tractable in this case. Recall from Eq. (9), setting $\epsilon = 0$, the transition distribution conditioned on $\boldsymbol{\theta}$ is a Gaussian whose mean is linear in $\boldsymbol{\theta}$.

Therefore the stationary distribution:

$$\mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{P}} \left[\mathbb{P}_{\mathbf{S}, \mathbb{P}^{\text{data}}} (\boldsymbol{\theta}' | \boldsymbol{\theta}) \right],$$

will be a Gaussian distribution as a sum of gaussians is also a gaussian. Let us assume the distribution has mean \mathbf{m} . Comparing the means we have:

$$\begin{aligned} \mathbf{m}(1 - \eta) + \eta\boldsymbol{\mu} &= \mathbf{m} \\ \implies \mathbf{m} &= \boldsymbol{\mu}. \end{aligned}$$

Moreover, $\mathbb{P}_{\mathbf{S}, \mathbb{P}^{\text{data}}} (\boldsymbol{\theta}' | \boldsymbol{\theta})$ is a Gaussian with covariance $\eta^2 \mathbf{S}$ for all $\boldsymbol{\theta}$. Hence the expectation over \mathbb{P} also has covariance $\eta^2 \mathbf{S}$. This concludes the proof. \square

Lemma B.2. *The loss at stationarity of the learning dynamics in the absence of an adversary for the mean estimation problem for $\mathbb{P}^{\text{data}} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is given by:*

$$\mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{P}(\mathbf{S}, \mathbb{P}^{\text{data}})} [\ell(\boldsymbol{\theta})] = \eta^2 \text{Trace}(\mathbf{S}). \quad (11)$$

Proof.

$$\begin{aligned} & \mathbb{E}_{\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}, \eta^2 \mathbf{S})} \left[\|\boldsymbol{\theta} - \boldsymbol{\mu}\|_2^2 \right] \\ &= \mathbb{E}_{\boldsymbol{\theta} \sim \mathcal{N}(0, \eta^2 \mathbf{S})} \left[\|\boldsymbol{\theta}\|_2^2 \right] \\ &= \eta^2 \text{Trace}(\mathbf{S}). \end{aligned}$$

\square

Remark B.1. *We use CVXPY [6] to solve the optimization problems in Algorithm 1.*