

---

# Certifying robustness to adaptive data poisoning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       The rise of foundational models fine-tuned with human feedback from  
2       potentially untrusted users has increased the risk of adversarial data poi-  
3       soning, necessitating the study of robustness of learning algorithms against  
4       such attacks. While existing research focuses on certifying robustness for  
5       static adversaries acting on offline datasets, dynamic attack algorithms have  
6       shown to be more effective. Relevant for models with periodic updates  
7       where an adversary can adapt based on the algorithm’s behavior, such  
8       as those in RLHF, we present a novel framework for computing certified  
9       bounds on the impact of dynamic poisoning, and use these certificates to  
10      design robust learning algorithms. We give an illustration of the framework  
11      for the mean-estimation problem.

## 12 1 Introduction & Problem Formulation

13      With the advent of foundational models fine tuned using human feedback gathered from  
14      potentially untrusted users (for example, users of a publicly available language model)  
15      [3, 7], the potential for adversarial or malicious data entering the training data of a model  
16      increases substantially . This motivates the study of robustness of learning algorithms to  
17      poisoning attacks [2]. More recently, there have been works that attempt to achieve “certified  
18      robustness” to data poisoning, i.e., proving that the worst case impact of poisoning is below  
19      a certain bound that depends on parameters of the learning algorithm. All the work in this  
20      space, to the best of our knowledge, focuses on the *static* poisoning adversary [9, 13]. Even  
21      in [10] which is the closest setting to our work, the poisoning adversary acts over offline  
22      datasets in a temporally extended fashion which are poisoned in one shot, and thus is not  
23      dynamic. There has been work on *dynamic* attack algorithms [12, 11] showing that these  
24      attacks can indeed be more powerful than static adversaries. This motivates the question  
25      we study: can we obtain certificates of robustness for a broad class of learning algorithms  
26      against *dynamic* poisoning adversaries?

27      In this paper, we study learning algorithms corrupted by a dynamic poisoning adversary  
28      who can observe the behavior of the learning algorithm and adapt the poisoning in response.  
29      This is relevant in scenarios where models are continuously / periodically updated in the face  
30      of new feedback, as is common in RLHF / fine tuning applications. We provide (to the best  
31      of our knowledge) the first general framework for computing certified bounds on the worst  
32      case impact of a data poisoning attacker, and further, use this certificate to design robust  
33      learning algorithms. We given an illustration of the framework for the mean-estimation  
34      problem (see Section 2), and aim to leverage this framework for regression, classification  
35      and generative models in future work.

36 **Learning objective** We study learning problems where the goal is to minimize

$$\mathbb{E}_{z \sim \mathbb{P}^{\text{data}}} [\ell(\boldsymbol{\theta}, z)]$$

37 where  $\boldsymbol{\theta} \in \mathbb{R}^d$  are parameters to be estimated (for example parameters of a generative model,  
38 classification model or a regression model),  $\ell : \mathbb{R}^d \times \mathbb{R}^n \mapsto \mathbb{R}$  is a loss function and  $z \in \mathbb{R}^n$   
39 are i.i.d. samples from an underlying data distribution  $\mathbb{P}^{\text{data}}$ .

40 **Adversarially corrupted learning algorithm** We work in a setting where the learning is  
41 being done in an online fashion and the corrupted datapoint can be updated after every  
42 step of learning, based on the trajectory of the learning process observed by the adversary.  
43 We consider learning algorithms of the form

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \eta (\nabla \ell(\boldsymbol{\theta}_t, z_t) + \mathbf{B}w_t) \quad (1)$$

44 where  $w_t \sim \mathcal{N}(0, \mathbf{I})$  is chosen iid at each  $t$  and  $z_t \sim \epsilon \text{Dirac}(z_t^{\text{adv}}) + (1 - \epsilon) \mathbb{P}^{\text{data}}$  and  
45  $\mathbf{B} \in \mathbb{R}^{d \times d}$  is a design parameter of the learning algorithm that is described below (see  
46 Potential Defense) and  $\epsilon$  is a parameter that controls the “level” of poisoning (analogous to  
47 the fraction of poisoned samples). This is a special case of Huber’s contamination model,  
48 which is used in the robust statistics literature [4] (with the contamination model being a  
49 Dirac distribution). Further, there are typically allowed ranges for the datapoints that come  
50 from the learning algorithm normalizing inputs or by an outlier detection system used to  
51 filter potential adversarial data. In this preliminary work, we restrict ourselves to norm  
52 balls,  $\mathcal{A} = \{z : \|z\| \leq r\}$ .

53 **Potential Defense** Inspired by differentially private learning algorithms like DP-SGD  
54 [1], we propose adding Gaussian noise to the learning process as a way of smoothing the  
55 learning algorithm against impacts of the poisoning adversary. In particular, we add  $\mathbf{B}w_t$   
56 where  $w_t$  is iid noise in each step sampled from the standard Gaussian, and  $\mathbf{B}$  is a design  
57 parameter of the learning algorithm. Subsequently, we will choose  $\mathbf{B}$  so as to minimize the  
58 worst case impact of the poisoning adversary. We denote by  $\mathbf{S} = \mathbf{B}\mathbf{B}^\top$  the covariance matrix  
59 of the noise added.

60 **Adversarial objective** We assume that the poisoning adversary is interested in maximiz-  
61 ing some adversarial objective  $\ell_{\text{adv}}(\boldsymbol{\theta})$  on target data:

$$\ell_{\text{adv}}(\boldsymbol{\theta}) = \mathbb{E}_{z \sim \mathbb{P}^{\text{target}}} [\ell(\boldsymbol{\theta}, z)] \text{ Maximize loss on some target data}$$

62 **Dynamics as a Markov Chain** By (1), we have that, conditioned on  $\boldsymbol{\theta}_t$  and  $z_t$ ,  $\boldsymbol{\theta}_{t+1}$  follows  
63 a Gaussian distribution with mean  $\boldsymbol{\theta}_t - \eta \nabla \ell(\boldsymbol{\theta}_t, z_t)$ .

64 The dynamics (1) gives rise to a Markov chain over the parameters  $\boldsymbol{\theta}$ . If  $\mathbb{P}_t$  denotes the  
65 distribution over parameters at time  $t$ , we have

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

66 where  $\mathbb{P}_{\mathbf{S}, \mathbb{P}^{\text{data}}, z^{\text{adv}}}$  is the transition kernel induced by (1), explicitly given by

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

67 where  $\mathcal{N}(x | \mu, \Sigma)$  denotes the Gaussian density at  $x$  for a Gaussian with mean  $\mu$  and covari-  
68 ance matrix  $\Sigma$ .

69 **A certificate for the adversarial loss (Analysis)** Since this is a Markov process, the optimal  
70 sequence of actions for the adversary (ie choices of  $z^{\text{adv}}$ ) constitute a Markov Decision  
71 Process with

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

72 and hence, can be formulated as an infinite dimensional linear program [8]. In particular,  
 73 for the infinite horizon average reward setting [6], the LP can be written as

$$\sup_{\mathbb{P} \in \mathcal{P}[\mathbb{R}^d \times \mathbb{R}^n]} \mathbb{E}_{\boldsymbol{\theta}, z^{\text{adv}} \sim \mathbb{P}} [\ell_{\text{adv}}(\boldsymbol{\theta})] \quad (3a)$$

$$\text{subject to } \mathbb{E}_{\boldsymbol{\theta}, z^{\text{adv}} \sim \mathbb{P}} \left[ \mathbb{P}_{S, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) \right] = \mathbb{E}_{\boldsymbol{\theta}, z^{\text{adv}} \sim \mathbb{P}} [\mathbb{I}[\boldsymbol{\theta}' = \boldsymbol{\theta}]] \quad \forall \boldsymbol{\theta}' \in \mathbb{R}^d. \quad (3b)$$

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

76 **Theorem 1.** For any function  $\lambda : \mathbb{R}^d \mapsto \mathbb{R}$ , we can upper bound the optimal value of (3) by

$$\sup_{\substack{\boldsymbol{\theta} \in \mathbb{R}^d \\ z^{\text{adv}} \in \mathcal{A}}} \mathbb{E}_{\boldsymbol{\theta}' \sim \mathbb{P}_{S, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\cdot | \boldsymbol{\theta})} [\lambda(\boldsymbol{\theta}')] + \ell_{\text{adv}}(\boldsymbol{\theta}) - \lambda(\boldsymbol{\theta}). \quad (4)$$

77 *Proof.* Follows by weak duality for the LP (5a).

78 If strong duality holds, we further have that the optimal value of (3) is exactly equal to

$$\inf_{\lambda: \mathbb{R}^d \mapsto \mathbb{R}} \sup_{\boldsymbol{\theta}, z^{\text{adv}}} \mathbb{E}_{\boldsymbol{\theta}' \sim \mathbb{P}_{S, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\cdot | \boldsymbol{\theta})} [\lambda(\boldsymbol{\theta}')] + \ell_{\text{adv}}(\boldsymbol{\theta}) - \lambda(\boldsymbol{\theta}). \quad (5a)$$

79 **A design principle for robust learning algorithms (aka meta-learning a robust learning**  
 80 **algorithm)** Based on the above analysis, we can attempt to design the parameters of the  
 81 learning algorithm (in this case  $S = \mathbf{B}\mathbf{B}^\top$ ) to trade-off performance and robustness. In  
 82 particular, in the absence of poisoned data, the updates (1) result in a stationary distribution  
 83  $\mathbb{P}(S, \mathbb{P}^{\text{data}})$  over model parameters  $\boldsymbol{\theta}$ :

$$\mathbb{P}(S, \mathbb{P}^{\text{data}}) = \mathbb{P} \text{ that satisfies } \mathbb{P}(\boldsymbol{\theta}') = \mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{P}} \left[ \mathbb{E}_{z \sim \mathbb{P}^{\text{data}}} \left[ \mathcal{N}(\boldsymbol{\theta}' | \boldsymbol{\theta}_t - \eta \nabla \ell(\boldsymbol{\theta}, z), \eta^2 S) \right] \right]. \quad (6)$$

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

$$\inf_{\substack{S \in \mathcal{S}_+^d \\ \lambda: \mathbb{R}^d \mapsto \mathbb{R}}} \mathbb{E}_{\mathbb{P}^{\text{data}} \sim \mathcal{P}} \left[ \mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{P}(S, \mathbb{P}^{\text{data}})} [\ell(\boldsymbol{\theta})] + \kappa \left( \sup_{\boldsymbol{\theta} \in \mathbb{R}^d, z^{\text{adv}} \in \mathcal{A}} \mathbb{E}_{\boldsymbol{\theta}' \sim \mathbb{P}_{S, \mathbb{P}^{\text{data}}, z^{\text{adv}}}(\cdot | \boldsymbol{\theta})} [\lambda(\boldsymbol{\theta}')] + \ell_{\text{adv}}(\boldsymbol{\theta}) - \lambda(\boldsymbol{\theta}) \right) \right], \quad (7)$$

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

## 92 2 Mean estimation

93 Consider the mean estimation problem, where we aim to learn the parameter  $\boldsymbol{\theta}$  to estimate  
 94 the mean  $\boldsymbol{\mu} = \mathbb{E}_{z \sim \mathbb{P}^{\text{data}}} [z]$  of a distribution  $\mathbb{P}^{\text{data}}$ . The adversarial loss is given by:

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

95 **Certificate on adversarial loss (analysis)**

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

$$\inf_{A \in \mathcal{S}^d, \mathbf{b} \in \mathbb{R}^d, \nu \geq 0} g(\mathbf{A}, \mathbf{b}, \nu, S, \boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (8)$$

---

**Algorithm 1** Meta learning
 

---

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

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

$$g(\mathbf{A}, \mathbf{b}, v, \mathbf{S}, \mu, \Sigma) = \begin{cases} \frac{1}{4} \left\| \begin{bmatrix} 2(1-\epsilon)\eta(1-\eta)\mathbf{A}\mu - 2\mu - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} + 2v\mu \end{bmatrix} \right\|_D^2 + (1-\epsilon)(\eta^2 \operatorname{Trace}(\Sigma\mathbf{A}) + \eta^2 \mu^\top \mathbf{A} \mu + \eta \mathbf{b}^\top \mu) \\ + \mu^\top \mu + \eta^2 \operatorname{Trace}(\mathbf{A}\mathbf{S}) + v(r - \mu^\top \mu) \text{ if } v \geq 0; \mathbf{D} \succeq 0 \\ -\infty \text{ else} \end{cases} \quad (9)$$

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

102 **Meta-Learning Algorithm** Following the formulation in Eq. (7), we wish to learn a defense  
 103 parameter  $\mathbf{S}$  that minimizes the expected loss (expectation over different  $\mathbb{P}^{\text{data}}$  from the  
 104 meta distribution  $\mathcal{P}$ ). For the mean estimation problem this boils down to solving:

$$\inf_{\mathbf{S} \in \mathbb{S}_+^d} \eta^2 \operatorname{Trace}(\mathbf{S}) + \kappa \mathbb{E}_{\mu, \Sigma \sim \mathcal{P}} \left[ \inf_{\substack{v \geq 0 \\ \mathbf{A} \in \mathbb{S}^d, \mathbf{b} \in \mathbb{R}^d}} g(\mathbf{A}, \mathbf{b}, v, \mathbf{S}, \mu, \Sigma) \right]. \quad (10)$$

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

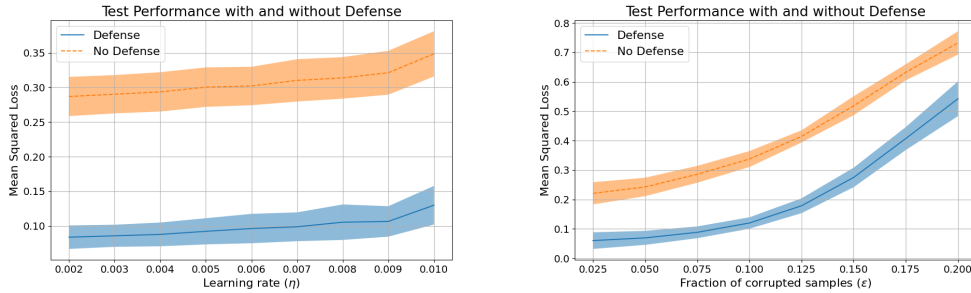


Figure 1: Test performance (mean squared error between true and estimated means) on 50 ( $d = 20$  dimensional) Gaussian distributions drawn from Gaussian prior for the mean and Inverse Wishart prior for the covariance. The defense parameter  $\mathbf{S}$  was trained with 10 such randomly chosen Gaussians 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 significantly.

## References

- 108
- 109 [1] Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimiza-  
110 tion: Efficient algorithms and tight error bounds. In *2014 IEEE 55th annual symposium*  
111 *on foundations of computer science*, pages 464–473. IEEE, 2014.
- 112 [2] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support  
113 vector machines. *arXiv preprint arXiv:1206.6389*, 2012.
- 114 [3] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei.  
115 Deep reinforcement learning from human preferences. *Advances in neural information*  
116 *processing systems*, 30, 2017.
- 117 [4] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*.  
118 Cambridge university press, 2023.
- 119 [5] Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language  
120 for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- 121 [6] Alan Malek, Yasin Abbasi-Yadkori, and Peter Bartlett. Linear programming for large-  
122 scale markov decision problems. In *International conference on machine learning*, pages  
123 496–504. PMLR, 2014.
- 124 [7] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela  
125 Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training  
126 language models to follow instructions with human feedback. *Advances in neural*  
127 *information processing systems*, 35:27730–27744, 2022.
- 128 [8] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*.  
129 John Wiley & Sons, 2014.
- 130 [9] Jacob Steinhardt, Pang Wei W Koh, and Percy S Liang. Certified defenses for data  
131 poisoning attacks. *Advances in neural information processing systems*, 30, 2017.
- 132 [10] Wenxiao Wang and Soheil Feizi. Temporal robustness against data poisoning. *Advances*  
133 *in Neural Information Processing Systems*, 36, 2024.
- 134 [11] Yizhen Wang and Kamalika Chaudhuri. Data poisoning attacks against online learning.  
135 *arXiv preprint arXiv:1808.08994*, 2018.
- 136 [12] Xuezhou Zhang, Xiaojin Zhu, and Laurent Lessard. Online data poisoning attacks. In  
137 *Learning for Dynamics and Control*, pages 201–210. PMLR, 2020.
- 138 [13] Yuhao Zhang, Aws Albarghouthi, and Loris D’Antoni. Bagflip: A certified defense  
139 against data poisoning. *Advances in Neural Information Processing Systems*, 35:31474–  
140 31483, 2022.

141 **A Proofs**

142 **Theorem 2.** Choosing  $\lambda : \mathbb{R}^d \rightarrow \mathbb{R}$  in Theorem 1 to be quadratic, i.e.  $\lambda(\boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{A} \boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{b}$ ,  
 143 the adversarial constraint set of the form  $\|\mathbf{z}^{\text{adv}} - \boldsymbol{\mu}\|_2^2 \leq r$ , the certificate for the mean estimation  
 144 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:

$$\inf_{\mathbf{A} \in \mathbb{S}^d, \mathbf{b} \in \mathbb{R}^d, \nu \geq 0} g(\mathbf{A}, \mathbf{b}, \nu, \mathbf{S}, \boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (8)$$

145 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  
 146 Matrix Inequality (LMI) constraint) as defined below:

$$g(\mathbf{A}, \mathbf{b}, \nu, \mathbf{S}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \begin{cases} \frac{1}{4} \left\| \begin{bmatrix} 2(1-\epsilon)\eta(1-\eta)\mathbf{A}\boldsymbol{\mu} - 2\boldsymbol{\mu} - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} + 2\nu\boldsymbol{\mu} \end{bmatrix} \right\|_{\mathbf{D}}^2 + (1-\epsilon)(\eta^2 \text{Trace}(\boldsymbol{\Sigma}\mathbf{A}) + \eta^2 \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} + \eta \mathbf{b}^\top \boldsymbol{\mu}) \\ + \boldsymbol{\mu}^\top \boldsymbol{\mu} + \eta^2 \text{Trace}(\mathbf{A}\mathbf{S}) + \nu(r - \boldsymbol{\mu}^\top \boldsymbol{\mu}) \text{ if } \nu \geq 0; \mathbf{D} \succeq 0 \\ -\infty \text{ else} \end{cases} \quad (9)$$

147 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  $\|\mathbf{x}\|_{\mathbf{D}}^2 = \mathbf{x}^\top \mathbf{D}^{-1} \mathbf{x}$ .

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

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

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

152 The transition distribution for the parameter is given by:

$$\mathbb{P}_{\mathbf{S}, \mathbb{P}^{\text{data}}, \mathbf{z}^{\text{adv}}}(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \epsilon \mathcal{N}(\boldsymbol{\theta}' | F(\boldsymbol{\theta}, \mathbf{z}^{\text{adv}}), \eta^2 \mathbf{S}) + (1 - \epsilon) \mathbb{E}_{\mathbf{z} \sim \mathbb{P}^{\text{data}}} \left[ \mathcal{N}(\boldsymbol{\theta}' | F(\boldsymbol{\theta}, \mathbf{z}), \eta^2 \mathbf{S}) \right] \quad (11)$$

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

154 Then, we have from Eq. (5a) that the certified bound on the adversarial objective is given by:  
 155

$$\sup_{\theta} \epsilon \mathbb{E}_{\mathbf{z}^{\text{adv}} \in \mathcal{A}} \mathbb{E}_{\theta' \sim \mathcal{N}(F(\theta, \mathbf{z}^{\text{adv}}), \eta^2 \mathbf{S})} [\lambda(\theta')] + (1 - \epsilon) \mathbb{E}_{\mathbf{z} \sim \mathbb{P}^{\text{data}}} \left[ \mathbb{E}_{\theta' \sim \mathcal{N}(F(\theta, \mathbf{z}), \eta^2 \mathbf{S})} [\lambda(\theta')] \right] - \lambda(\theta) + \ell_{\text{adv}}(\theta) \quad (12a)$$

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

$$= \sup_{\theta} \epsilon \left( \lambda(F(\theta, \mathbf{z}^{\text{adv}})) + \eta^2 \langle \nabla^2 \lambda(0), \mathbf{S} \rangle \right) + (1 - \epsilon) \mathbb{E}_{\mathbf{z} \sim \mathbb{P}^{\text{data}}} \left[ \lambda(F(\theta, \mathbf{z})) + \eta^2 \langle \nabla^2 \lambda(0), \mathbf{S} \rangle \right] - \lambda(\theta) + \ell_{\text{adv}}(\theta) \quad (12b)$$

$$= \sup_{\theta} \left[ - \left\| \begin{bmatrix} \theta \\ \mathbf{z}^{\text{adv}} \end{bmatrix} \right\|_{\mathbf{E}^{-1}}^2 + \begin{bmatrix} \theta \\ \mathbf{z}^{\text{adv}} \end{bmatrix}^\top \begin{bmatrix} 2(1 - \epsilon)\eta(1 - \eta)\mathbf{A}\mu - 2\mu - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} \end{bmatrix} \right] + (1 - \epsilon)(\eta^2 \text{Trace}(\mathbf{\Sigma}\mathbf{A}) + \eta^2 \mu^\top \mathbf{A}\mu + \eta\mathbf{b}^\top \mu) + \mu^\top \mu, \quad (12c)$$

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

The dual function of this supremum (with dual variable  $\nu$ ) can be written as:

$$= \inf_{\nu \geq 0} \sup_{\theta} \left[ - \left\| \begin{bmatrix} \theta \\ \mathbf{z}^{\text{adv}} \end{bmatrix} \right\|_{\mathbf{D}^{-1}}^2 + \begin{bmatrix} \theta \\ \mathbf{z}^{\text{adv}} \end{bmatrix}^\top \begin{bmatrix} 2(1 - \epsilon)\eta(1 - \eta)\mathbf{A}\mu - 2\mu - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} \end{bmatrix} \right] + (1 - \epsilon)(\eta^2 \text{Trace}(\mathbf{\Sigma}\mathbf{A}) + \eta^2 \mu^\top \mathbf{A}\mu + \eta\mathbf{b}^\top \mu) + \mu^\top \mu + \nu(r - \mu^\top \mu) \quad (12d)$$

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

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

$$\inf_{\nu \geq 0} \frac{1}{4} \left\| \begin{bmatrix} 2(1 - \epsilon)\eta(1 - \eta)\mathbf{A}\mu - 2\mu - \eta\mathbf{b} \\ \epsilon\eta\mathbf{b} + 2\nu\mu \end{bmatrix} \right\|_{\mathbf{D}}^2 + (1 - \epsilon)(\eta^2 \text{Trace}(\mathbf{\Sigma}\mathbf{A}) + \eta^2 \mu^\top \mathbf{A}\mu + \eta\mathbf{b}^\top \mu) + \mu^\top \mu + \eta^2 \text{Trace}(\mathbf{A}\mathbf{S}) + \nu(r - \mu^\top \mu) \text{ such that } \mathbf{D} \succeq 0. \quad (12e)$$

156 This completes the proof.

157

□

158 **Lemma A.1.** *The stationary distribution in the absence of adversary in Eq. (6) for the mean*  
 159 *estimation problem for  $\mathbb{P}^{\text{data}} = \mathcal{N}(\mu, \mathbf{\Sigma})$  takes the form:*

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

160 *Proof.* The stationary distribution is tractable in this case. Recall from Eq. (11), setting  
 161  $\epsilon = 0$ , the transition distribution conditioned on  $\theta$  is a Gaussian whose mean is linear in  $\theta$ .  
 162 Therefore the stationary distribution:

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

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

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

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

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

$$\mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{P}(\mathcal{S}, \mathbb{P}^{\text{data}})}[\ell(\boldsymbol{\theta})] = \eta^2 \text{Trace}(\mathbf{S}). \quad (13)$$

*Proof.*

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

169

$\square$

170 **Remark A.1.** *We use CVXPY [5] to solve the optimization problems in Algorithm 1.*