

---

# Certified defences hurt generalisation

---

Piersilvio De Bartolomeis  
ETH Zürich

Jacob Clarysse  
ETH Zürich

Fanny Yang  
ETH Zürich

Amartya Sanyal  
ETH Zürich

## Abstract

In recent years, much work has been devoted to designing certified defences for neural networks, i.e., methods for learning neural networks that are provably robust to certain adversarial perturbations. Due to the non-convexity of the problem, dominant approaches in this area rely on convex approximations, which are inherently loose. In this paper, we question the effectiveness of such approaches for realistic computer vision tasks. First, we provide extensive empirical evidence to show that certified defences suffer not only worse accuracy but also worse robustness and fairness than empirical defences. We hypothesise that the reason for why certified defences suffer in generalisation is (i) the large number of relaxed non-convex constraints and (ii) the strong alignment between the adversarial perturbations and the "signal" direction. We provide a combination of theoretical and experimental evidence to support these hypotheses.

## 1 Introduction

Deep neural networks have been shown to be vulnerable to adversarial attacks: imperceptible perturbations to the input can fool state-of-the-art classifiers [13, 27, 21]. The existence of adversarial examples raises serious security concerns in many safety-critical applications [29, 10], and so robustness to adversarial attacks is becoming a crucial design goal for modern deep learning architectures.

In practice, robustness against many different types of perturbations may be desired depending on the domain of application. Hence, in order to build robust models, we need to first define a threat model for the adversary, i.e., a set of perturbations  $\mathcal{B}$ . The most commonly used threat models in the literature are norm-bounded perturbations, where  $\mathcal{B}_{\epsilon,p} = \{\delta : \|\delta\|_p \leq \epsilon\}$  is the  $\ell_p$ -ball with radius  $\epsilon$  centred around the origin. In this paper, we focus on  $\mathcal{B}_{\epsilon,2}$  which, for ease of notation, we will represent as  $\mathcal{B}_\epsilon$ .

For any distribution  $\mathcal{D}$ , neural network model  $f_\theta$  parameterised by  $\theta$ , and loss function  $L$ , our goal is to find a model which solves the following robust optimisation problem:

$$\min_{\theta} \mathbf{R}_\epsilon(\theta) := \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\delta \in \mathcal{B}_\epsilon} L(f_\theta(x + \delta), y) \right] \quad (1)$$

We call  $\mathbf{R}_\epsilon(\theta)$  the robust error when  $L$  is the 0-1 loss function. In practice, as the distribution  $\mathcal{D}$  is unknown, we instead minimise the empirical robust error on a finite dataset  $D$  sampled from  $\mathcal{D}$ . Unfortunately, in the case of neural networks, the inner maximisation is a non-convex optimisation problem and prohibitively hard to solve from a computational perspective [16, 30]. Therefore, two kinds of approaches are widely used to efficiently solve it: *empirical* defences that provide a lower bound on the solution and *certified* defences that provide an upper bound.

Among empirical defences, Adversarial Training (AT) [13, 20] is one of the few that has stood the test of time. AT minimises the worst-case empirical loss in Equation (1) by approximately solving the inner-maximisation problem with first-order gradient-based optimisation methods. However, despite its simplicity and computational efficiency, owing to its heuristic nature, AT does not provide any robustness guarantee. In many safety critical domains, such guarantees are of immense importance.

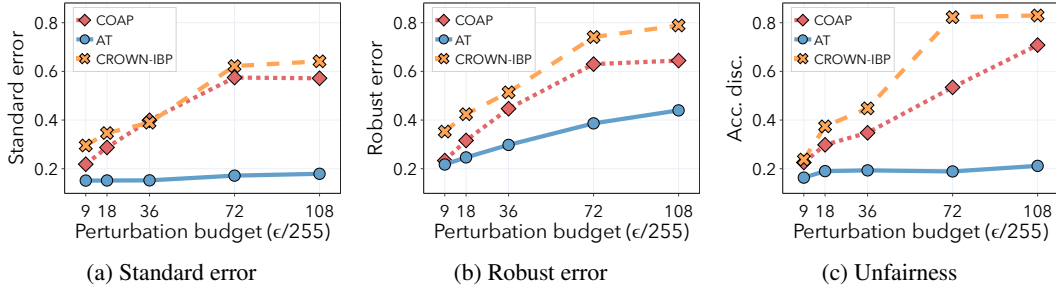


Figure 1: Results for  $\ell_2$ -adversaries on the CIFAR-10 dataset. We compare ResNet architectures trained using state-of-the-art certified defenses CROWN-IBP [37, 34] and COAP [32, 31] against the most popular empirical defense to date AT [20, 13]. In Figures 1a, 1b and 1c we plot respectively standard error, robust error and accuracy discrepancy as the perturbation budget increases. See Appendix D.3 for complete experimental details.

To address this limitation, recently, there has been significant interest in designing certified defences, i.e., methods for learning neural networks that are *provably* robust to norm-bounded perturbations on the training data. Many recent works [31, 23, 7, 37] have proposed to solve a convex relaxation of the inner-maximisation problem by relaxing the non-convex ReLU constraint sets with convex ones. Despite all of these progresses, certified defences based on convex relaxations suffer from an inherent flaw: the upper bound they provide on the robust error is far from being tight due to the looseness of the convex relaxation [25]. In this paper, we argue that the fundamental looseness of convex relaxations hinders the practical effectiveness of current certified defences. In particular, as shown in Figure 1, certified defences suffer significantly worse accuracy, robustness, and fairness on the test data compared to adversarial training. Our contributions are as follows:

- In Section 2, we show that current certified defences hurt accuracy, robustness, and fairness across a range of  $\ell_2$ -ball perturbations on real-world vision datasets like MNIST and CIFAR-10.
- In Section 3, we provide experimental evidence that certified defences hurt generalisation because of (i) the large number of relaxed non-convex constraints and (ii) the strong alignment between the adversarial perturbations and the signal direction.
- In Section 4, we prove, in a simplified high-dimensional classification setting, that certified defences yield higher robust risk than adversarial training when the adversarial perturbation aligns with the signal direction.

## 2 Certified defences hurt generalisation on real-world data

In this section, we show that certified defences hurt standard error, robust error, and fairness on two common computer vision datasets: MNIST [19] and CIFAR-10 [18]. Among certified defences, we consider the convex outer adversarial polytope (COAP) approach [32, 31], which achieves state-of-the-art certified robustness under  $\ell_2$ -ball perturbations. Additionally, we consider CROWN-IBP [36, 34], which uses the tightest convex relaxation CROWN [36] and achieves state-of-the-art certified robustness under  $\ell_\infty$ -ball perturbations. Among empirical defences, we consider adversarial training (AT) [20, 13], which is one of the most popular and effective defences to date.

**Models and robust evaluation** We consider the  $\ell_2$ -ball perturbations threat model. To reliably evaluate the robust error, we use the strongest version of AutoAttack (AA+) [4]. For CIFAR-10, we train a residual network (ResNet) and for MNIST we train a large convolutional neural network (CNN). Both architectures were introduced in Wong et al. [32] as standard benchmarks for certified defences. We refer the reader to Appendix D.3 for experimental details.

**Certified defences hurt standard and robust error** Several studies have shown that adversarial training may lead to an increase in standard error when compared with standard training [24, 28, 35]. Here, we observe the same phenomenon to a much higher degree in certified defences. Specifically, our experimental results show that certified defences not only suffer worse standard error but also worse robust error than adversarial training. First, we observe on both MNIST and CIFAR-10 in Figures 2a and 2c, respectively, that for increasing perturbation budget, the standard error gap between

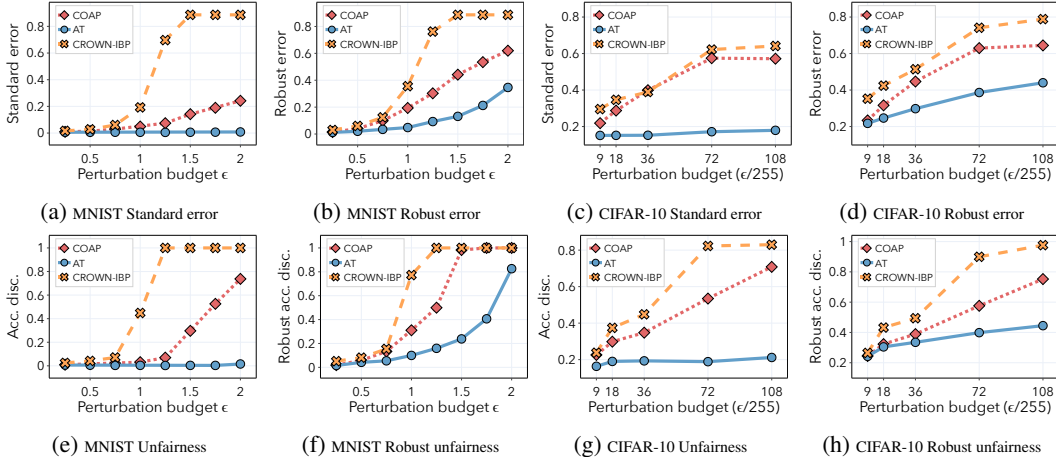


Figure 2: Results for  $\ell_2$ -adversaries on MNIST and CIFAR-10 datasets. In Figures 2a, 2b, 2e and 2f we plot respectively the standard error, robust error, accuracy discrepancy and robust accuracy discrepancy for a CNN trained on MNIST, as the perturbation budget  $\epsilon$  increases. In Figures 2c, 2d, 2g and 2h we plot respectively the standard error, robust error, accuracy discrepancy and robust accuracy discrepancy, for a ResNet trained on CIFAR-10, as the perturbation budget  $\epsilon$  increases.

certified (CROWN-IBP, COAP) and empirical defences (AT) increases. In particular, the gap reaches almost 90% for CROWN-IBP on MNIST when  $\epsilon = 1.5$ . Secondly, we observe that the robust error gap increases with increasing perturbation budget for both MNIST and CIFAR-10 in Figures 2b and 2d, respectively. In particular, the gap reaches almost 40% for the largest perturbation budgets.

**Certified defences hurt fairness** Previously, we showed that certified defences hurt both standard and robust generalisation. Taking it one step further, we show that certified defences (CROWN-IBP, COAP) suffer significantly worse fairness than empirical defences (AT).

Let  $\mathbf{R}(\theta)$  be the standard error of the classifier  $f_\theta$  and  $\mathbf{R}^k(\theta)$  the standard error conditioned on the class label  $k$ . We measure the degree of unfairness as:

$$\frac{\max_k \mathbf{R}^k(\theta) - \mathbf{R}(\theta)}{1 - \mathbf{R}(\theta)} \quad (2)$$

Using the terminology in Sanyal et al. [26], we refer to this metric as *accuracy discrepancy*. However, we expect our results to translate to other related fairness metrics as well [6, 8, 14, 15]. Additionally, we consider the discrepancy in robust accuracy, as it was observed in Xu et al. [33] that adversarial defences may induce a large discrepancy of robustness among different classes. We refer to this metric as *robust accuracy discrepancy* and it corresponds to replacing standard error with the robust error in Equation (2).

We present our experimental results comparing the fairness of certified and empirical defences. For MNIST, we observe in Figure 2e and 2f that COAP and CROWN-IBP have a significant discrepancy for both standard and robust accuracy. For large perturbations, these methods obtain 100% discrepancy, indicating that their accuracy on the worst class can be as low as 0%. By contrast, AT preserves fairness for both standard and robust accuracy much better. In particular, the discrepancy for standard accuracy is always less than 2% for all perturbation budgets considered. Similarly, for CIFAR-10 AT maintains a constant accuracy discrepancy around 20% for all perturbation budgets considered, whereas for certified defences it steadily increases with the perturbation budget, reaching above 80%. Additionally, for robust accuracy, we observe a discrepancy gap of 35% between the best certified and empirical defences for the largest perturbation budget considered.

### 3 Developing intuition on synthetic datasets

In this section, we hypothesise that certified defences hurt robust and standard generalisation because of (i) the large number of relaxed non-convex constraints and (ii) the strong alignment between the

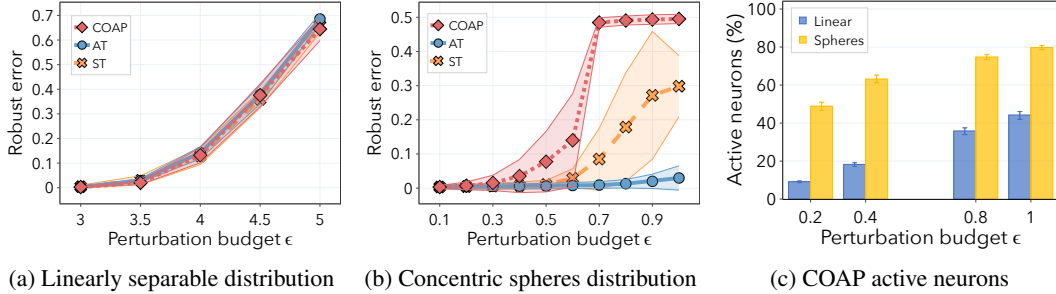


Figure 3: We report mean and standard deviation over 15 seeds. In Figures 3a and 3b we plot the robust error for standard training (ST), adversarial training (AT) and convex outer adversarial polytope (COAP), when training on the linearly separable and concentric spheres distributions respectively. In Figure 3c, we plot the percentage of neurons in the activation set for the linearly separable and concentric spheres distribution respectively. See Appendix D.2 for complete experimental details.

adversarial perturbations and the signal direction. We investigate these hypotheses on more controlled settings. In particular, we consider two synthetic data distributions: a linearly separable distribution as in Clarysse et al. [3], which is similar to distributions studied in Nagarajan and Kolter [22], Tsipras et al. [28], and the concentric spheres distribution studied in Gilmer et al. [12], Nagarajan and Kolter [22].

**Data and threat models** Similarly to the previous section, we focus on  $\ell_2$ -ball perturbations of size  $\epsilon$ . As for distributions, we consider the linearly separable distribution where first, the label  $y \in \{+1, -1\}$  is drawn with equal probability. Then, for some  $\gamma > 0$ , the covariate vector is created as  $x = [\gamma \text{sgn}(y); \tilde{x}]$ , where  $\tilde{x} \in \mathbb{R}^{d-1}$  is a random vector drawn from a standard normal distribution  $\tilde{x} \sim \mathcal{N}(0, \sigma^2 I_{d-1})$  and  $[\cdot; \cdot]$  represents concatenation. We sample the concentric spheres dataset as follows; for  $0 < R_1 < R_{-1}$ , we first draw a binary label  $y \in \{+1, -1\}$  with equal probability and then the covariate vector  $x \in \mathbb{R}^d$  is distributed uniformly on the sphere of radius  $R_y$ . Observe that achieving a low test error on the concentric spheres distribution requires a non-linear classifier.

In Figures 3a and 3b, we plot the robust error of standard training (ST), adversarial training (AT), and certified training (COAP) on the linear and concentric spheres distributions respectively. We see that in contrast to the linear setting, COAP has a much higher robust error on the concentric spheres distribution than AT and ST, where the gap increases with increasing perturbation budget  $\epsilon$ . Below, we provide intuition as to why COAP has a much higher robust error than AT on the concentric spheres distribution.

The intuition is two-fold: first of all COAP relaxes the non-convex ReLU constraints for all neurons that activate within the perturbation set, i.e., there exists  $\delta \in \mathcal{B}_\epsilon$  for which the input to the neuron equals 0. Hence, the larger the percentage of relaxed neurons, the worse the approximation. This is formally captured by Theorem A.2 in Appendix A. Secondly, the  $\ell_2$ -ball perturbations are significantly aligned with the signal direction, meaning that they effectively reduce the information about the label in the data. Applying an approximation in this direction yields poor generalisation. We prove this in Theorem 4.1 for the linearly separable distribution.

**COAP relaxes many constraints on the concentric spheres** In Figure 3c we empirically show that COAP convexly approximates a large number of constraints when training on the concentric spheres distribution. We plot the percentage of active neurons on the concentric spheres and linear distributions against increasing perturbation budgets: the percentage is much higher for the concentric spheres than for the linearly separable distribution and increases with perturbation budget  $\epsilon$ . Indeed, the complex spherical decision boundary requires much more active neurons compared to the linear decision boundary which only needs 1 active neuron.

**$\ell_2$ -ball perturbations align with the signal direction** We empirically show that  $\ell_2$ -ball perturbations align with the signal direction on the concentric spheres distribution. Note that for a point  $x$  drawn from the concentric spheres distribution, the signal direction is given by  $y \frac{x}{\|x\|_2}$  (see Figure 4a for a 2D visualization). In Figure 4b, we plot the cosine distance between the  $\ell_2$ -perturbations

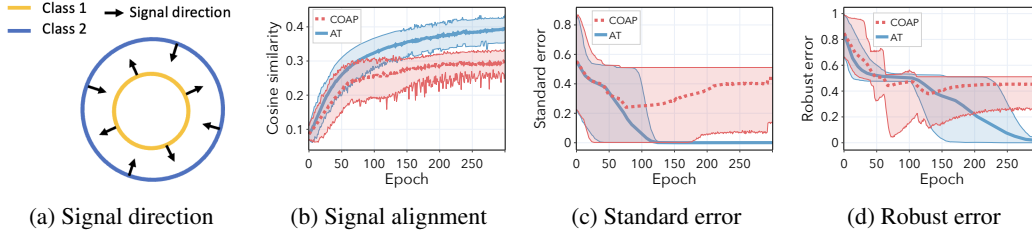


Figure 4: We report mean and standard deviation over 15 seeds. In Figure 4a we plot a 2-dimensional realisation of the adversarial spheres dataset, the black arrow illustrates the signal direction. In Figure 4b we plot the cosine similarity between  $\ell_2$  norm-bounded perturbations on the training data (average) and the signal directed vector. In Figures 4c and 4d we plot standard and robust error for adversarial training (AT) and convex outer adversarial polytope (COAP). We observe that when cosine similarity is high, the gap in standard and robust error between COAP and AT increases. Hence, especially approximations in the signal direction can hurt standard and robust generalisation.

computed on the training set, and the signal direction. Comparing Figures 4b to 4d, we see that during the early stages of training, the  $\ell_2$ -ball perturbations are not aligned with the signal direction and the robust and standard errors for COAP are similar to AT. However, after some epochs, when the  $\ell_2$ -ball perturbations start to align with the signal direction, both the robust and standard error gaps between COAP and AT increase. This provides evidence that, as training progresses,  $\ell_2$ -ball perturbations become significantly aligned with the signal direction and the generalisation gap worsens.

#### 4 Approximations along the signal direction hurt generalisation

In this section, we further investigate our hypothesis that certified defences hurt generalisation when adversarial perturbations are aligned with the signal direction. In particular, we study the linearly separable distribution from the previous section and assume that the adversarial attacks concentrate all of their perturbation budget along the signal direction. In Theorem 4.1, we prove for a simple neural network that, in high dimensions, certified defences (COAP) yield higher robust error than empirical defences (AT) for large perturbation budgets. We then corroborate our theoretical results with extensive experimental evidence on synthetic data.

**Data and threat models** We consider the linearly separable distribution described in Section 3. As for the threat model, we consider signal-directed attacks that efficiently concentrate their attack budget on the signal in the input. Since the signal direction corresponds to the first component of the data, we define the set of allowed perturbations as:

$$\mathcal{B}_\epsilon(x) = \{z_1 = x + e_1\beta \mid |\beta| \leq \epsilon\} \quad (3)$$

where  $e_1$  is the standard basis vector of the first coordinate. Further, as the original formulation of COAP only allows  $\ell_p$ -adversaries, we provide an extension of COAP that covers our theoretical and experimental setting in Appendix A.

**One gradient step training** We consider the hypothesis class to be the set of one-neuron shallow neural networks  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ , defined by:

$$f_\theta(x) = a \text{ReLU}(\theta^\top x) + b \quad (4)$$

where  $x \in \mathbb{R}^d, \theta \in \mathbb{R}^d, a \in \mathbb{R}, b \in \mathbb{R}$  and the only trainable parameter is  $\theta_1$ . Note that as our distribution is linearly separable, our hypothesis class includes the ground truth.

We study the early phase of neural network optimisation. Under structural assumptions on the data, it has been proved that one gradient step with sufficiently large learning rate can drastically decrease the training loss [2] and extract task-relevant features [11, 5]. A similar setting was also studied recently in Ba et al. [1] for the MSE loss in the high-dimensional asymptotic limit. Here, we focus on the classification setting with binary cross-entropy loss. Below we state our main theorem.

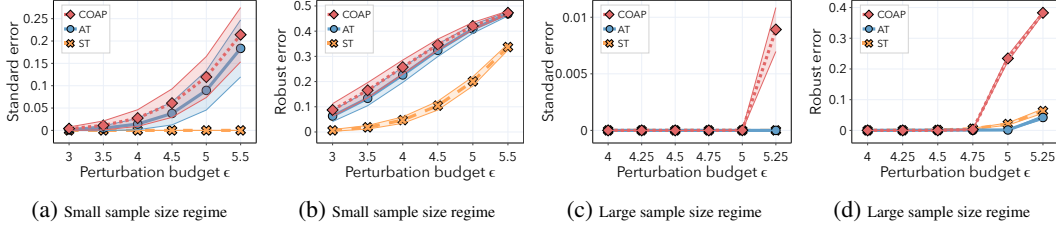


Figure 5: We report mean and standard deviation over 15 seeds. In Figure 5a and 5b we plot respectively the standard and robust errors in the small sample size ( $n = 50$ ) regime for standard training (ST), adversarial training (AT) and convex outer adversarial polytope (COAP) as the perturbation budget  $\epsilon$  increases. In Figure 5c and 5d we plot respectively the standard and robust errors in the large sample size ( $n = 10000$ ) regime for standard training (ST), adversarial training (AT) and convex outer adversarial polytope (COAP) as the perturbation budget  $\epsilon$  increases. See Appendix D.1 for complete experimental details.

**Theorem 4.1.** Let  $\bar{\theta}$  and  $\tilde{\theta}$  be the network parameters after one step of gradient descent with respect to AT and COAP objectives. Let,

$$\frac{\|\theta_{2:d}\|_2}{\|\theta_1\|_2} > \sqrt{\frac{24\gamma^3}{\sigma^2}} \quad \text{and} \quad \frac{2}{3}\gamma < \epsilon < \gamma \quad (5)$$

where  $\theta$  are the network parameters at initialization. Then, COAP yields higher robust risk than AT:

$$\mathbf{R}_\epsilon(\tilde{\theta}) > \mathbf{R}_\epsilon(\bar{\theta}) \quad (6)$$

Theorem 4.1 relies on two main assumptions. The first is an assumption on the data dimensionality and the initialisation of the network parameters  $\theta$ . For instance, if the network parameters are initialised sampling from a standard multivariate gaussian  $\theta \sim \mathcal{N}(0, I_d)$ , then we know that  $\|\theta\|_2$  concentrates around  $\sqrt{d}$  with high probability. Hence, the assumption is satisfied when the data dimensionality  $d$  is sufficiently high. Further, the second assumption requires that the perturbation budget  $\epsilon$  is sufficiently close to the separation margin  $\gamma$ . This is consistent with the experimental evidence we presented so far, as the generalisation of certified defences significantly worsen for large perturbation budgets.

**Synthetic experiments** We corroborate our theory with experimental evidence using a one-hidden layer neural network with 100 neurons. In particular, we investigate the effect of perturbation budget  $\epsilon$  on generalisation for three different models: standard training (ST), adversarial training (AT) [20, 13] and convex outer adversarial polytope (COAP) [31, 32]. In Figure 5, we plot robust and standard errors for both small and large sample size regimes as the perturbation budget  $\epsilon$  increases. The generalisation gap in the small sample size regime between standard and adversarial training was already observed in Clarysse et al. [3] for linear classifiers. Here, we observe a further generalisation gap between AT and COAP in both small and large sample size regimes, which surprisingly worsens in the large sample regime.

## 5 Conclusions

In this paper, we show that certified defences can hurt accuracy, robustness and fairness for realistic datasets and adversarial perturbations. Further, we develop intuition on synthetic datasets for why certified defences hurt generalisation, combining both theoretical and experimental evidence. We believe that shedding light on the performance gap between empirical and certified defences will not only provide us with a clearer picture of the trade-offs observed in practice but also lead to better approaches for adversarial robustness.

## References

- [1] Jimmy Ba, Murat A. Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation, May 2022. arXiv:2205.01445 [cs, math, stat].

- [2] Niladri S. Chatterji, Philip M. Long, and Peter L. Bartlett. When Does Gradient Descent with Logistic Loss Find Interpolating Two-Layer Networks? *Journal of Machine Learning Research*, (159), 2021. ISSN 1533-7928.
- [3] Jacob Clarysse, Julia Hörrmann, and Fanny Yang. Why adversarial training can hurt robust accuracy, 2022. arXiv:2203.02006.
- [4] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *Proceedings of the International Conference on Machine Learning*, 2020.
- [5] Amit Daniely and Eran Malach. Learning Parities with Neural Networks. In *Advances in Neural Information Processing Systems*, 2020.
- [6] John C. Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization. *CoRR*, abs/1810.08750, 2018. URL <http://arxiv.org/abs/1810.08750>.
- [7] Krishnamurthy Dvijotham, Robert Stanforth, Sven Gowal, Timothy A. Mann, and Pushmeet Kohli. A Dual Approach to Scalable Verification of Deep Networks. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, 2018.
- [8] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard S. Zemel. Fairness through awareness. In Shafi Goldwasser, editor, *Innovations in Theoretical Computer Science 2012, Cambridge, MA, USA, January 8-10, 2012*, pages 214–226. ACM, 2012. doi: 10.1145/2090236.2090255. URL <https://doi.org/10.1145/2090236.2090255>.
- [9] Ecenaz Erdemir, Jeffrey Bickford, Luca Melis, and Sergül Aydıre. Adversarial Robustness with Non-uniform Perturbations. In *Advances in Neural Information Processing Systems*, 2021.
- [10] Samuel G. Finlayson, John D. Bowers, Joichi Ito, Jonathan L. Zittrain, Andrew L. Beam, and Isaac S. Kohane. Adversarial attacks on medical machine learning. *Science*, (6433), March 2019.
- [11] Spencer Frei, Niladri S. Chatterji, and Peter L. Bartlett. Random Feature Amplification: Feature Learning and Generalization in Neural Networks, May 2022. arXiv:2202.07626 [cs, math, stat].
- [12] Justin Gilmer, Luke Metz, Fartash Faghri, Samuel S. Schoenholz, Maithra Raghu, Martin Wattenberg, and Ian J. Goodfellow. Adversarial Spheres. *CoRR*, 2018. arXiv: 1801.02774.
- [13] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and Harnessing Adversarial Examples. In *Proceedings of the International Conference on Learning Representations*, 2015.
- [14] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 3315–3323, 2016. URL <https://proceedings.neurips.cc/paper/2016/hash/9d2682367c3935defcb1f9e247a97c0d-Abstract.html>.
- [15] Úrsula Hébert-Johnson, Michael P. Kim, Omer Reingold, and Guy N. Rothblum. Multicalibration: Calibration for the (computationally-identifiable) masses. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1944–1953. PMLR, 2018. URL <http://proceedings.mlr.press/v80/hebert-johnson18a.html>.
- [16] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks. In *Proceedings of the International Conference of Computer Aided Verification*, 2017.
- [17] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *Proceedings of the International Conference on Learning Representations*, 2015.

- [18] Alex Krizhevsky. Learning multiple layers of features from tiny images. *citeseer*, 2009.
- [19] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, (11), 1998.
- [20] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards Deep Learning Models Resistant to Adversarial Attacks. In *Proceedings of the International Conference on Learning Representations*, 2018.
- [21] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. DeepFool: A Simple and Accurate Method to Fool Deep Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [22] Vaishnavh Nagarajan and J. Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. In *Advances in Neural Information Processing Systems*, 2019.
- [23] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified Defenses against Adversarial Examples. In *Proceedings of the International Conference on Learning Representations*, 2018.
- [24] Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John C. Duchi, and Percy Liang. Understanding and Mitigating the Tradeoff between Robustness and Accuracy. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 7909–7919. PMLR, 2020. URL <http://proceedings.mlr.press/v119/raghunathan20a.html>.
- [25] Hadi Salman, Greg Yang, Huan Zhang, Cho-Jui Hsieh, and Pengchuan Zhang. A Convex Relaxation Barrier to Tight Robustness Verification of Neural Networks. In *Advances in Neural Information Processing Systems*, 2019.
- [26] Amartya Sanyal, Yaxi Hu, and Fanny Yang. How unfair is private learning? In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, 2022.
- [27] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *Proceedings of the International Conference on Learning Representations*, 2014.
- [28] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness May Be at Odds with Accuracy. In *Proceedings of the International Conference on Learning Representations*, 2019.
- [29] Cumhur Erkan Tuncali, Georgios Fainekos, Hisahiro Ito, and James Kapinski. Simulation-based Adversarial Test Generation for Autonomous Vehicles with Machine Learning Components. In *Proceedings of the IEEE Intelligent Vehicles Symposium*, 2018.
- [30] Tsui-Wei Weng, Huan Zhang, Hongge Chen, Zhao Song, Cho-Jui Hsieh, Luca Daniel, Duane S. Boning, and Inderjit S. Dhillon. Towards fast computation of certified robustness for relu networks. In *Proceedings of the International Conference on Machine Learning*, 2018.
- [31] Eric Wong and J. Zico Kolter. Provable Defenses against Adversarial Examples via the Convex Outer Adversarial Polytope. In *Proceedings of the International Conference on Machine Learning*, 2018.
- [32] Eric Wong, Frank R. Schmidt, Jan Hendrik Metzen, and J. Zico Kolter. Scaling provable adversarial defenses. In *Advances in Neural Information Processing Systems*, 2018.
- [33] Han Xu, Xiaorui Liu, Yaxin Li, Anil K. Jain, and Jiliang Tang. To be Robust or to be Fair: Towards Fairness in Adversarial Training. In *Proceedings of the International Conference on Machine Learning*, Proceedings of Machine Learning Research, 2021.
- [34] Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya Kailkhura, Xue Lin, and Cho-Jui Hsieh. Automatic Perturbation Analysis for Scalable Certified Robustness and Beyond. In *Advances in Neural Information Processing Systems*, 2020.



- [35] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I. Jordan. Theoretically Principled Trade-off between Robustness and Accuracy. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 7472–7482. PMLR, 2019. URL <http://proceedings.mlr.press/v97/zhang19p.html>.
- [36] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient Neural Network Robustness Certification with General Activation Functions. In *Advances in Neural Information Processing Systems*, 2018.
- [37] Huan Zhang, Hongge Chen, Chaowei Xiao, Sven Gowal, Robert Stanforth, Bo Li, Duane S. Boning, and Cho-Jui Hsieh. Towards Stable and Efficient Training of Verifiably Robust Neural Networks. In *Proceedings of the International Conference on Learning Representations*, 2020.

## A Certified defences for signal-directed adversaries

We now formulate the convex outer adversarial polytope (COAP) [31] for adversaries that concentrate all their budget along the signal direction in the input. Our derivation can be seen as an extension of Wong and Kolter [31], Erdemir et al. [9].

**Data and threat model** We consider the linearly separable distribution described in Section 3. First, the label  $y \in \{+1, -1\}$  is drawn with equal probability. Then, for some  $\gamma > 0$ , the covariate vector is  $x = [\gamma \operatorname{sgn}(y); \tilde{x}]$ , where  $\tilde{x} \in \mathbb{R}^{d-1}$  is a random vector drawn from a standard normal distribution  $\tilde{x} \sim \mathcal{N}(0, \sigma^2 I_{d-1})$  and  $[\cdot; \cdot]$  represents concatenation. As for the threat model, we consider signal-directed attacks that efficiently concentrate their attack budget on the signal in the input. Since the signal direction corresponds to the first component of the data, we define the set of allowed perturbations as:

$$\mathcal{B}_\epsilon(x) = \{z_1 = x + e_1\beta \mid |\beta| \leq \epsilon\} \quad (7)$$

where  $e_1$  is the standard basis vector of the first coordinate.

**Network structure** For the rest of this section we consider a 2 layer feed-forward ReLU network. In particular, we define  $f_\theta(x) : \mathbb{R}^d \rightarrow \mathbb{R}^2$  as follows:

$$x \xrightarrow{x+\delta} z_1 \xrightarrow{W_1 z_1 + b_1} \hat{z}_2 \xrightarrow{\operatorname{ReLU}(\cdot)} z_2 \xrightarrow{W_2 z_2 + b_2} \hat{z}_3 \quad (8)$$

where  $x \in \mathbb{R}^d$ ,  $z_1 \in \mathcal{B}_\epsilon(x)$ ,  $W_1$  and  $W_2$ , are linear operators and we define the set of network parameters as  $\theta = \{W_i, b_i\}_{i=1,2}$ .

**Constructing the convex outer bound** We define the adversarial polytope  $\mathcal{Z}_\epsilon(x)$  as the set of all final-layer activations attainable by perturbing  $x$  with some  $\tilde{x} \in \mathcal{B}_\epsilon(x)$ :

$$\mathcal{Z}_\epsilon(x) = \{f_\theta(\tilde{x}) : \tilde{x} \in \mathcal{B}_\epsilon(x)\} \quad (9)$$

Our approach will be to construct a convex outer bound on this adversarial polytope: if no adversarial example exists in this outer approximation, then we are guaranteed that no adversarial example exists in the original polytope. We relax the ReLU activations  $z = \operatorname{ReLU}(\hat{z})$  with their convex envelopes:

$$z \geq 0, \quad z \geq \hat{z}, \quad (u - \ell)z \leq u\hat{z} - u\ell \quad (10)$$

where  $u$  and  $\ell$  are respectively the pre-activations  $\hat{z}$  upper and lower bounds.

First, we address the problem of obtaining the upper and lower bounds  $u$  and  $\ell$  for the pre-activations  $\hat{z}$ . Specifically, the following proposition gives a closed form solution for  $\ell$  and  $u$ .

**Proposition A.1.** *Consider the neural network  $f_\theta$  defined in Equation (8). Let  $w_1$  be the first column of  $W_1$ . Then, for a data point  $x$  and perturbation budget  $\epsilon$ , we have the following element-wise bounds on the pre-activation vector  $\hat{z}_2$ :*

$$\ell \leq \hat{z}_2 \leq u \quad (11)$$

where

$$\ell = W_1 x + b_1 - \epsilon |w_1| \quad \text{and} \quad u = W_1 x + b_1 + \epsilon |w_1| \quad (12)$$

*Proof.* Given a data point  $x$  and perturbation budget  $\epsilon$ , let  $\tilde{x} = x + \delta$  be the perturbed input to the network. First, we find an upper bound the pre-activations values  $\hat{z}_2$ :

$$\hat{z}_2 = W_1(x + \delta) + b_1 = W_1 x + b_1 + W_1 \delta \quad (13)$$

In particular, we want to solve the following optimisation problem for each component of the pre-activation vector:

$$u_i = \max_{\tilde{x} \in \mathcal{B}_\epsilon(x)} [\hat{z}_2]_i = [W_1 x]_i + [b_1]_i + \max_{\tilde{x} \in \mathcal{B}_\epsilon(x)} [W_1 \delta]_i \quad (14)$$

where  $u$  will be the vector containing element-wise upper bounds. Note that  $\delta = \beta e_1$ , thus the optimisation problem can be rewritten as:

$$\max_{\tilde{x} \in \mathcal{B}_\epsilon(x)} [W_1 \delta]_i = \max_{\|\beta\|_1 \leq \epsilon} \beta \cdot [w_1]_i = \epsilon \cdot \|[w_1]_i\|_1 \quad (15)$$

where  $w_1$  is the first column of  $W_1$ . The vector of upper bounds will then be:

$$u = W_1 x + b_1 + \epsilon |w_1| \quad (16)$$

Along the same lines, we can derive the vector of lower bounds  $\ell$ :

$$l = W_1 x + b_1 - \epsilon |w_1| \quad (17)$$

□

Next, we define the outer bound on the adversarial polytope we get from relaxing ReLU constraints as  $\tilde{\mathcal{Z}}_\epsilon(x)$ . Given a data point  $x$  with known label  $y$ , we can formulate the problem of finding an adversarial example with a linear program as follows:

$$\min_{\hat{z}_3} [\hat{z}_3]_y - [\hat{z}_3]_{\bar{y}} = c^\top \hat{z}_3 \quad \text{s.t. } \hat{z}_3 \in \tilde{\mathcal{Z}}_\epsilon(x) \quad (18)$$

where  $\bar{y}$  is the binary negation of  $y$ . Note that if we solve this linear program and find that the objective is positive, then we know that no input perturbation within the threat model can misclassify the example.

However, solving the linear program in Equation (18) for every example in the dataset is intractable. Therefore, we consider the dual formulation and take a feasible solution. In the following theorem, we state the dual problem formulation of the linear program in Equation (18).

**Theorem A.2.** *The dual of the linear program (18) can be written as*

$$\begin{aligned} \max_{\alpha} \quad & \tilde{J}_\epsilon(x, g_\theta(c, \alpha)) \\ \text{s.t.} \quad & \alpha_j \in [0, 1], \forall j \end{aligned} \quad (19)$$

where  $\tilde{J}_\epsilon(x, \nu_1, \nu_2, \nu_3)$  is equal to

$$-\sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|[\hat{\nu}_1]_1\|_1 \quad (20)$$

and  $g_\theta$  is a one-hidden layer neural network given by the equations

$$\begin{aligned} \nu_3 &= -c \\ \hat{\nu}_2 &= W_2^\top \nu_3 \\ [\nu_2]_j &= 0, \quad j \in \mathcal{I}^- \\ [\nu_2]_j &= [\hat{\nu}_2]_j, \quad j \in \mathcal{I}^+ \\ [\nu_2]_j &= \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \alpha_j [\hat{\nu}_2]_j^-, \quad j \in \mathcal{I} \\ \hat{\nu}_1 &= W_1^\top \nu_2 \end{aligned} \quad (21)$$

where  $\mathcal{I}^-$ ,  $\mathcal{I}^+$  and  $\mathcal{I}$  denote the sets of activations in the hidden layer where  $\ell$  and  $u$  are both negative, both positive or span zero, respectively.

*Proof.* Consider a data point  $x$  and let  $\tilde{x} = x + \delta$  be the adversarial perturbed data point. First, we explicit all the constraints for the linear program defined in (18):

$$\begin{aligned} \min_{\hat{z}_3} \quad & [\hat{z}_3]_y - [\hat{z}_3]_{\bar{y}} = c^\top \hat{z}_3, \quad \text{s.t.} \\ & x + \delta \in \mathcal{B}_\epsilon(x) \\ & z_1 = x + \delta \\ & \hat{z}_2 = W_1 z_1 + b_1 \\ & \hat{z}_3 = W_2 z_2 + b_2 \\ & [z_2]_j = 0, \quad \forall j \in \mathcal{I}^- \\ & [z_2]_j = [\hat{z}_2]_j, \quad \forall j \in \mathcal{I}^+ \\ & [z_2]_j \geq 0, \quad \forall j \in \mathcal{I} \\ & [z_2]_j \geq [\hat{z}_2]_j, \quad \forall j \in \mathcal{I} \\ & ((u_j - \ell_j) [z_2]_j - u_j [\hat{z}_2]_j) \leq -u_j \ell_j, \quad \forall j \in \mathcal{I} \end{aligned} \quad (22)$$

where  $\mathcal{I}^-$ ,  $\mathcal{I}^+$  and  $\mathcal{I}$  denote the sets of activations in the hidden layer where  $\ell$  and  $u$  are both negative, both positive, or span zero respectively. In order to compute the dual of this problem, we associate the following Lagrangian variables with each of the constraints:

$$\begin{aligned}
\hat{z}_2 &= W_1 z_1 + b_1 \Rightarrow \nu_2 \\
\hat{z}_3 &= W_2 z_2 + b_2 \Rightarrow \nu_3 \\
z_1 &= x + \delta \Rightarrow \psi \\
-[z_2]_j &\leq 0 \Rightarrow \mu_j, \forall j \in \mathcal{I} \\
[\hat{z}_2]_j - [z_2]_j &\leq 0 \Rightarrow \tau_j, \forall j \in \mathcal{I} \\
((u_j - \ell_j) [z_2]_j - u_j [\hat{z}_2]_j) &\leq -u_j \ell_j \Rightarrow \lambda_j, \forall j \in \mathcal{I}
\end{aligned} \tag{23}$$

note that we do not define explicit dual variables for  $[z_2]_j = 0$  and  $[z_2]_j = [\hat{z}_2]_j$  as we can easily eliminate them. We write the Lagrangian as follows:

$$\begin{aligned}
\mathcal{L}(z, \hat{z}, \nu, \delta, \lambda, \tau, \mu, \psi) &= - (W_1^\top \nu_2 + \psi)^\top z_1 - \sum_{j \in \mathcal{I}} \left( \mu_j + \tau_j - \lambda_j (u_j - \ell_j) + [W_2^\top \nu_3]_j \right) [z_2]_j \\
&+ \sum_{j \in \mathcal{I}} (\tau_j - \lambda_j u_j + [\nu_2]_j) [\hat{z}_2]_j + (c + \nu_3)^\top \hat{z}_3 - \sum_{i=1}^2 \nu_{i+1}^\top b_i \\
&+ \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x + \psi^\top \delta + \sum_{j \in \mathcal{I}^-} [\hat{z}_2]_j [\nu_2]_j \\
&+ \sum_{j \in \mathcal{I}^+} [z_2]_j ([\nu_2]_j - [W_2^\top \nu_3]_j) \\
\text{s.t. } \tilde{x} &\in \mathcal{B}_\epsilon(x)
\end{aligned} \tag{24}$$

and we take the infimum w.r.t.  $z, \hat{z}, \delta$ :

$$\begin{aligned}
\inf_{z, \hat{z}, \delta} \mathcal{L}(z, \hat{z}, \nu, \delta, \lambda, \tau, \mu, \psi) &= - \inf_{z_2} \sum_{j \in \mathcal{I}} \left( \mu_j + \tau_j - \lambda_j (u_j - \ell_j) + [W_2^\top \nu_3]_j \right) [z_2]_j \\
&+ \inf_{\hat{z}_2} \sum_{j \in \mathcal{I}} (\tau_j - \lambda_j u_j + [\nu_2]_j) [\hat{z}_2]_j + \inf_{\hat{z}_3} (c + \nu_3)^\top \hat{z}_3 - \sum_{i=1}^2 \nu_{i+1}^\top b_i \\
&+ \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x + \inf_{\tilde{x} \in \mathcal{B}_\epsilon(x)} \psi^\top \delta - \inf_{z_1} (W_1^\top \nu_2 + \psi)^\top z_1 \\
&+ \inf_{\hat{z}_2} \sum_{j \in \mathcal{I}^-} [\hat{z}_2]_j [\nu_2]_j + \inf_{z_2} \sum_{j \in \mathcal{I}^+} [z_2]_j ([\nu_2]_j - [W_2^\top \nu_3]_j)
\end{aligned} \tag{25}$$

Now, we can compute the infimum for the  $\psi^\top \delta$  term:

$$\inf_{\tilde{x} \in \mathcal{B}_\epsilon(x)} \psi^\top \delta = \inf_{\|\beta\|_1 \leq \epsilon} \psi_1 \cdot \beta = -\epsilon \cdot \|\psi_1\|_1 \tag{26}$$

and since for all the other terms the infimum of a linear function is  $-\infty$ , except in the special case when it is identically zero, the infimum of  $\mathcal{L}(\cdot)$  becomes:

$$\inf_{z, \hat{z}, \delta} \mathcal{L}(\cdot) = \begin{cases} -\sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x - \epsilon \|\psi_1\|_1 & \text{if conditions} \\ -\infty & \text{else} \end{cases} \tag{27}$$

where the conditions to satisfy are:

$$\begin{aligned}
\nu_3 &= -c \\
W_1^\top \nu_2 &= -\psi \\
[\nu_2]_j &= 0, j \in \mathcal{I}_i^- \\
[\nu_2]_j &= [W_2^\top \nu_3]_j, j \in \mathcal{I}_i^+ \\
\left. \begin{aligned} (u_j - \ell_j) \lambda_j - \mu_j - \tau_j &= [W_2^\top \nu_3]_j \\ [\nu_2]_j &= u_j \lambda_j - \tau_j \end{aligned} \right\} j \in \mathcal{I} \\
\lambda, \tau, \mu &\geq 0
\end{aligned} \tag{28}$$

Thus, we can rewrite the dual problem as follows:

$$\begin{aligned}
& \max_{\nu, \psi, \lambda, \tau, \mu} - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x - \epsilon \|\psi_1\|_1 \\
& \text{s.t.} \quad \nu_3 = -c \\
& W_1^\top \nu_2 = -\psi \\
& [\nu_2]_j = 0, j \in \mathcal{I}^- \\
& [\nu_2]_j = [W_2^\top \nu_3]_j, j \in \mathcal{I}^+ \\
& \left. \begin{aligned} (u_j - \ell_j) \lambda_j - \mu_j - \tau_j &= [W_2^\top \nu_3]_j \\ [\nu_2]_j &= u_j \lambda_j - \tau_j \end{aligned} \right\} j \in \mathcal{I} \\
& \lambda, \tau, \mu \geq 0
\end{aligned} \tag{29}$$

Note that the dual variable  $\lambda$  corresponds to the upper bounds in the convex ReLU relaxation, while  $\mu$  and  $\tau$  correspond to the lower bounds. By the complementarity property, we know that at the optimal solution, these variables will be zero if the ReLU constraint is non-tight, or non-zero if the ReLU constraint is tight. Since the upper and lower bounds cannot be tight simultaneously, either  $\lambda$  or  $\mu + \tau$  must be zero. This means that at the optimal solution to the dual problem we can decompose  $[W_2^\top \nu_3]_j$  into positive and negative parts since  $(u_j - \ell_j)\lambda_j \geq 0$  and  $\tau_j + \mu_j \geq 0$ :

$$\begin{aligned}
(u_j - \ell_j)\lambda_j &= [W_2^\top \nu_3]_j^+ \\
\tau_j + \mu_j &= [W_2^\top \nu_3]_j^-
\end{aligned} \tag{30}$$

combining this with the constraint  $[\nu_2]_j = u_j \lambda_j - \tau_j$  leads to

$$[\nu_2]_j = \frac{u_j}{u_j - \ell_j} [W_2^\top \nu_3]_j^+ - \alpha_j [W_2^\top \nu_3]_j^- \tag{31}$$

for  $j \in \mathcal{I}$  and  $0 \leq \alpha_j \leq 1$ . Hence, we have that:

$$\lambda_j = \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ \tag{32}$$

Now, we denote  $\hat{\nu}_1 = -\psi$  to make our notation consistent, and putting all of this together the dual objective becomes:

$$\begin{aligned}
- \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \lambda_j u_j \ell_j + \psi^\top x - \epsilon \|\psi_1\|_1 &= - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \frac{u_j \ell_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1 \\
&= - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1
\end{aligned} \tag{33}$$

and the final dual problem:

$$\begin{aligned}
& \max_{\nu, \hat{\nu}} - \sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1 \\
& \text{s.t.} \quad \nu_3 = -c \\
& \quad \hat{\nu}_2 = W_2^\top \nu_3 \\
& \quad [\nu_2]_j = 0, j \in \mathcal{I}^- \\
& \quad [\nu_2]_j = [\hat{\nu}_2]_j, j \in \mathcal{I}^+ \\
& \quad [\nu_2]_j = \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \alpha_j [\hat{\nu}_2]_j^-, j \in \mathcal{I} \\
& \quad \hat{\nu}_1 = W_1^\top \nu_2
\end{aligned} \tag{34}$$

□

Most importantly, the theorem states that we can represent the dual problem as a linear back propagation network, which provides a tractable solution for a lower bound on the primal objective. In practice, rather than solving the exact dual problem, we choose a fixed dual feasible solution:

$$\alpha_j = \frac{u_j}{u_j - \ell_j} \tag{35}$$

## B Theoretical results for signal-directed adversaries

In this section, we study a simplified one-neuron neural network. We consider the linearly separable distribution from the previous section and assume that the adversarial attacks concentrate all of their perturbation budget along the signal direction. For a formal definition of the data and threat models we refer the reader to Appendix A.

**One-neuron neural network** We consider the hypothesis class to be the set of one-neuron shallow neural networks  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ , defined by:

$$f_\theta(x) = a \operatorname{ReLU}(\theta^\top x) + b \quad (36)$$

where  $x \in \mathbb{R}^d, \theta \in \mathbb{R}^d, a \in \mathbb{R}, b \in \mathbb{R}$  and the only trainable parameter is  $\theta_1$ . Note that as our distribution is linearly separable, our hypothesis class includes the ground truth. Further, we focus on the classification setting with binary cross-entropy loss:

$$L(x, y) = y \log(\sigma(x)) + (1 - y) \log(1 - \sigma(x)) \quad (37)$$

where  $\sigma(\cdot)$  is the sigmoid function.

### B.1 Adversarial training

The basic idea behind adversarial training is to update the network parameters according to the following rule:

$$\theta \leftarrow \theta - \frac{\eta}{|D|} \sum_{(x, y) \in D} \nabla_\theta \max_{x + \delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x + \delta), y) \quad (38)$$

This is usually done by applying some first-order approximation to the maximisation problem. However, for our simplified network we can analytically compute the gradient.

First of all, note that when  $L$  is the binary cross-entropy loss function we can rewrite the maximisation problem as follows:

$$\max_{x + \delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x + \delta), y) = L \left( \operatorname{sgn}(y) \overbrace{\min_{x + \delta \in \mathcal{B}_\epsilon(x)} \operatorname{sgn}(y) f_\theta(x + \delta)}^{:= J_\epsilon(x, y)}, y \right) \quad (39)$$

In particular, if  $J_\epsilon$  is strictly positive then no adversarial example exists that fools the network. Further, note that this formulation is closely related to the objective considered in Appendix A for the convex outer adversarial polytope. This will be useful when comparing COAP and AT gradients. Below we provide the gradient of the adversarial training objective w.r.t. the network parameters  $\theta$ .

**Proposition B.1.** *Consider the neural network  $f_\theta$  defined in Equation (36) and the threat model  $\mathcal{B}_\epsilon$  defined in Equation (7). Let  $L$  be the binary cross-entropy loss function, as defined in Equation (37). Then, we have:*

$$\begin{aligned} & \nabla_{\theta_1} \max_{x + \delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x + \delta), y) \\ &= -\operatorname{sgn}(y) \sigma(-J_\epsilon(x, y)) \begin{cases} a(x_1 - \epsilon \operatorname{sgn}(\theta_1)) \mathbf{1}\{\ell > 0\} & \text{if } a \operatorname{sgn}(y) > 0 \\ a(x_1 + \epsilon \operatorname{sgn}(\theta_1)) \mathbf{1}\{u > 0\} & \text{if } a \operatorname{sgn}(y) < 0 \end{cases} \end{aligned}$$

where  $\ell = \theta^\top x - \epsilon|\theta_1|$  and  $u = \theta^\top x + \epsilon|\theta_1|$  are respectively lower and upper bounds on the ReLU inputs.

*Proof.* Given a data point  $x$  with known label  $y \in \{-1, 1\}$ , when  $L$  is the binary cross-entropy loss function we have:

$$\max_{x + \delta \in \mathcal{B}_\epsilon(x)} L(f_\theta(x + \delta), y) = L \left( \operatorname{sgn}(y) \min_{x + \delta \in \mathcal{B}_\epsilon(x)} \operatorname{sgn}(y) f_\theta(x + \delta), y \right) \quad (40)$$

For our simplified network we can analytically compute a closed form solution of the minimisation problem:

$$\begin{aligned}
J_\epsilon &:= \min_{x+\delta \in \mathcal{B}_\epsilon(x)} \text{sgn}(y) (b + a \text{ReLU}(\theta^\top(x + \delta))) \\
&= \begin{cases} \text{sgn}(y) (b + a \max(0, \ell)) & \text{if } a \text{sgn}(y) > 0 \\ \text{sgn}(y) (b + a \max(0, u)) & \text{if } a \text{sgn}(y) < 0 \end{cases} \\
&= \begin{cases} \text{sgn}(y) (b + a \max(0, \ell)) & \text{if } a \text{sgn}(y) > 0 \\ \text{sgn}(y) (b + a \max(0, u)) & \text{if } a \text{sgn}(y) < 0 \end{cases}
\end{aligned} \tag{41}$$

where  $\ell = \theta^\top x - \epsilon|\theta_1|$  and  $u = \theta^\top x + \epsilon|\theta_1|$  are respectively lower and upper bounds on the pre-activations. Thus, we can compute the gradients for adversarial training w.r.t the signal parameter:

$$\frac{\partial}{\partial \theta_1} J_\epsilon = \begin{cases} \text{sgn}(y) a (x_1 - \epsilon \text{sgn}(\theta_1)) \mathbf{1}\{\ell > 0\} & \text{if } a \text{sgn}(y) > 0 \\ \text{sgn}(y) a (x_1 + \epsilon \text{sgn}(\theta_1)) \mathbf{1}\{u > 0\} & \text{if } a \text{sgn}(y) < 0 \end{cases} \tag{42}$$

and applying the chain-rule we have:

$$\frac{\partial}{\partial \theta_1} L(\text{sgn}(y) J_\epsilon, y) \tag{43}$$

$$= \frac{\partial}{\partial J_\epsilon} L(\text{sgn}(y) J_\epsilon, y) \cdot \frac{\partial}{\partial \theta_1} J_\epsilon \tag{44}$$

$$= \text{sgn}(y) [\sigma(\text{sgn}(y) J_\epsilon) - \mathbf{1}\{y = 1\}] \cdot \frac{\partial}{\partial \theta_1} J_\epsilon \tag{45}$$

$$= -\text{sgn}(y) \sigma(-J_\epsilon) \begin{cases} a(x_1 - \epsilon \text{sgn}(\theta_1)) \mathbf{1}\{\ell > 0\} & \text{if } a \text{sgn}(y) > 0 \\ a(x_1 + \epsilon \text{sgn}(\theta_1)) \mathbf{1}\{u > 0\} & \text{if } a \text{sgn}(y) < 0 \end{cases} \tag{46}$$

where in the last equality we use a known property of the sigmoid function,  $\sigma(x) = 1 - \sigma(-x)$ .  $\square$

## B.2 Convex outer adversarial polytope

We now consider the dual approximation  $\tilde{J}_\epsilon$  to the optimisation problem in Equation (39). Note that, for a binary classification problem, we have  $c = \text{sgn}(y)$  and the dual objective in Theorem A.2 becomes:

$$\tilde{J}_\epsilon(x, g_\theta(c, \alpha)) = \tilde{J}_\epsilon(x, y) \tag{47}$$

where we set  $\alpha$  to the dual feasible solution and for the sake of clarity we omit the dependence on the network parameters  $\theta$ .

We are particularly interested in the data points for which  $J_\epsilon(x, y) \neq \tilde{J}_\epsilon(x, y)$ , i.e., when the certified and adversarial training objectives differ. Below, we provide a necessary and sufficient condition to have a mismatch between the two objectives.

**Proposition B.2.** *Consider the neural network  $f_\theta$  defined in Equation (36) and the threat model  $\mathcal{B}_\epsilon$  defined in Equation (7). Let  $L$  be the binary cross-entropy loss function, as defined in Equation (37). Further, we define  $\ell = \theta^\top x - \epsilon|\theta_1|$  and  $u = \theta^\top x + \epsilon|\theta_1|$  respectively as lower and upper bounds on the ReLU inputs.*

*Let  $\mathcal{I}^* = \{(x, y) : 0 \in [\ell, u] \wedge a \text{sgn}(y) > 0\}$ . Then, for data points in  $\mathcal{I}^*$ , we have that AT and COAP gradients differ:*

$$\nabla_{\theta_1} J_\epsilon(x, y) \neq \nabla_{\theta_1} \tilde{J}_\epsilon(x, y) \quad \forall (x, y) \in \mathcal{I}^*$$

and COAP gradient is given by:

$$\begin{aligned}
&\nabla_{\theta_1} L(\text{sgn}(y) \tilde{J}_\epsilon(x, y), y) \\
&= -\frac{a \text{sgn}(y) \sigma(-\tilde{J}_\epsilon(x, y))}{2\epsilon} \left( \frac{\ell}{\|\theta_1\|_1} (x_1 + \epsilon \text{sgn}(\theta_1)) + u \frac{x_1 \|\theta_1\|_1 - \theta^\top x \text{sgn}(\theta_1)}{\theta_1^2} \right)
\end{aligned}$$

Further, for data points that are not in  $\mathcal{I}^*$  we have that AT and COAP gradients are equivalent:

$$\nabla_{\theta_1} J_\epsilon(x, y) = \nabla_{\theta_1} \tilde{J}_\epsilon(x, y) \quad \forall (x, y) \notin \mathcal{I}^*$$

*Proof.* For the sake of clarity, we report here the definition of COAP objective from appendix A.

$$\tilde{J}_\epsilon(x, y) = -\sum_{i=1}^2 \nu_{i+1}^\top b_i + \sum_{j \in \mathcal{I}} \ell_j [\nu_2]_j^+ - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1 \quad (48)$$

Further, recall that the dual variables  $\nu$  are given by the following equations:

$$\begin{aligned} \nu_3 &= -c \\ \hat{\nu}_2 &= W_2^\top \nu_3 \\ [\nu_2]_j &= 0, \quad j \in \mathcal{I}^- \\ [\nu_2]_j &= [\hat{\nu}_2]_j, \quad j \in \mathcal{I}^+ \\ [\nu_2]_j &= \frac{u_j}{u_j - \ell_j} [\hat{\nu}_2]_j^+ - \alpha_j [\hat{\nu}_2]_j^-, \quad j \in \mathcal{I} \\ \hat{\nu}_1 &= W_1^\top \nu_2 \end{aligned} \quad (49)$$

where  $\mathcal{I}^-$ ,  $\mathcal{I}^+$  and  $\mathcal{I}$  denote the sets of activations in the hidden layer where  $\ell$  and  $u$  are both negative, both positive and span zero, respectively.

First, we consider the case when the neuron is always dead, i.e.,  $\ell < u < 0$ . The dual variables are:

$$\begin{aligned} \nu_3 &= -\text{sgn}(y) \\ \hat{\nu}_2 &= -a \text{sgn}(y) \\ \nu_2 &= 0 \\ \hat{\nu}_1 &= 0 \end{aligned} \quad (50)$$

Hence, AT and COAP objectives are equal in this case:

$$\tilde{J}_\epsilon = \text{sgn}(y)b = J_\epsilon \quad (51)$$

where the last equality follows from Equation (41).

Next, we consider the case when the neuron is always active, i.e.,  $0 < \ell < u$ . The dual variables are:

$$\begin{aligned} \nu_3 &= -\text{sgn}(y) \\ \hat{\nu}_2 &= -a \text{sgn}(y) \\ \nu_2 &= -a \text{sgn}(y) \\ \hat{\nu}_1 &= -a \text{sgn}(y) \cdot \theta \end{aligned} \quad (52)$$

and the dual objective becomes:

$$\tilde{J}_\epsilon = -\nu_3^\top b - \hat{\nu}_1^\top x - \epsilon \|\hat{\nu}_1\|_1 \quad (53)$$

$$= \text{sgn}(y) (b + a(\theta^\top x)) - \epsilon \|a \text{sgn}(y) \theta\| \quad (54)$$

$$= \begin{cases} \text{sgn}(y) (b + a\ell) & \text{if } a \text{sgn}(y) > 0 \\ \text{sgn}(y) (b + au) & \text{if } a \text{sgn}(y) < 0 \end{cases} \quad (55)$$

$$= J_\epsilon \quad (56)$$

where the last equality follows from the fact that  $0 < \ell < u$ .

Finally, we consider the case when the neuron is in the activation set  $\mathcal{I}$ , i.e.,  $\ell < 0 < u$ . The dual variables are:

$$\begin{aligned} \nu_3 &= -\text{sgn}(y) \\ \hat{\nu}_2 &= -a \text{sgn}(y) \\ \nu_2 &= -a \text{sgn}(y) \frac{u}{2\epsilon \|\theta_1\|_1} \\ \hat{\nu}_1 &= -a \text{sgn}(y) \frac{u}{2\epsilon \|\theta_1\|_1} \cdot \theta \end{aligned} \quad (57)$$

Here we have two cases, when  $\hat{\nu}_2 > 0$  we can rewrite the dual objective as:

$$\tilde{J}_\epsilon = \text{sgn}(y) (b + au) = J_\epsilon \quad (58)$$

and the two objectives coincide.



When  $\nu_2 < 0$  we can rewrite the dual objective as:

$$\tilde{J}_\epsilon = \text{sgn}(y) \left( b + \frac{au\ell}{2\epsilon \|\theta_1\|_1} \right) \neq J_\epsilon \quad (59)$$

Hence, the only case when COAP gradient differs from AT gradient is when  $\nu_2 < 0$  and the neuron belongs to the activation set  $\mathcal{I}$ .

We compute the partial derivative w.r.t. the signal parameter  $\theta_1$  in this case, by the chain rule we have:

$$\frac{\partial}{\partial \theta_1} L(\text{sgn}(y) \cdot \tilde{J}_\epsilon, y) \quad (60)$$

$$= \frac{\partial}{\partial \tilde{J}_\epsilon} L(\text{sgn}(y) \cdot \tilde{J}_\epsilon, y) \cdot \frac{\partial}{\partial \theta_1} \tilde{J}_\epsilon \quad (61)$$

$$= \text{sgn}(y) \left[ \sigma(\text{sgn}(y) \cdot \tilde{J}_\epsilon) - \mathbf{1}\{y = 1\} \right] \cdot \frac{\partial}{\partial \theta_1} \tilde{J}_\epsilon \quad (62)$$

$$= -\frac{a \text{sgn}(y) \sigma(-\tilde{J}_\epsilon)}{2\epsilon} \left( \frac{\ell}{\|\theta_1\|_1} (x_1 + \epsilon \text{sgn}(\theta_1)) + u \frac{x_1 \|\theta_1\|_1 - \theta^\top x \text{sgn}(\theta_1)}{\theta_1^2} \right) \quad (63)$$

□

### B.3 Signal-directed approximations hurt generalisation

In this section, we prove that convex relaxations along the signal direction hurt robust generalisation. First, in Lemma B.1 we relate the robust error of the classifier  $f_\theta$  to the  $\ell_2$ -norm of the signal parameter  $\theta_1$ . Specifically, we show that robust error monotonically decreases in  $\|\theta_1\|_2$ .

**Lemma B.1.** *Let  $f_\theta$  be the neural network defined in Equation (36) and  $\mathcal{B}_\epsilon$  the threat model defined in Equation (7). We define the robust risk  $\mathbf{R}_\epsilon$  of  $f_\theta$  as follows:*

$$\mathbf{R}_\epsilon(\theta) := \mathbb{P}_{(x,y)} [\exists z \in \mathcal{B}_\epsilon(x) : y \neq \text{sgn}(f_\theta(z))] \quad (64)$$

Then,  $\mathbf{R}_\epsilon(\theta)$  is monotonically decreasing in  $\|\theta_1\|_2$ .

*Proof.*

$$\begin{aligned} \mathbf{R}_\epsilon(\theta) &:= \mathbb{P}_{(x,y)} [\exists z \in \mathcal{B}_\epsilon(x) : y \neq \text{sgn}(f_\theta(z))] \\ &= \frac{1}{2} (\mathbb{P}_x [\theta^\top x < \|b\|_1 \mid x_1 = \gamma - \epsilon] + \mathbb{P}_x [\theta^\top x > \|b\|_1 \mid x_1 = \epsilon - \gamma]) \\ &= \frac{1}{2} \left( \mathbb{P}_x \left[ \sum_{i=2}^d x_i \theta_i < -\theta_1(\gamma - \epsilon) + \|b\|_1 \right] + \mathbb{P}_x \left[ \sum_{i=2}^d x_i \theta_i > \theta_1(\gamma - \epsilon) + \|b\|_1 \right] \right) \\ &= \frac{1}{2} \left( \Phi \left( -\frac{(\gamma - \epsilon) \|\theta_1\|_2}{\sigma \|\theta_{2:d}\|_2} + \frac{\|b\|_1}{\sigma \|\theta_{2:d}\|_2} \right) + \Phi \left( -\frac{(\gamma - \epsilon) \|\theta_1\|_2}{\sigma \|\theta_{2:d}\|_2} - \frac{\|b\|_1}{\sigma \|\theta_{2:d}\|_2} \right) \right) \end{aligned}$$

hence  $\mathbf{R}_\epsilon(\theta)$  is monotonically decreasing in  $\|\theta_1\|_2$  and the statement follows. □

Below we present our main result.

**Theorem 4.1.** *Let  $\bar{\theta}$  and  $\tilde{\theta}$  be the network parameters after one step of gradient descent with respect to AT and COAP objectives. Let,*

$$\frac{\|\theta_{2:d}\|_2}{\|\theta_1\|_2} > \sqrt{\frac{24\gamma^3}{\sigma^2}} \quad \text{and} \quad \frac{2}{3}\gamma < \epsilon < \gamma \quad (5)$$

where  $\theta$  are the network parameters at initialization. Then, COAP yields higher robust risk than AT:

$$\mathbf{R}_\epsilon(\tilde{\theta}) > \mathbf{R}_\epsilon(\bar{\theta}) \quad (6)$$

*Proof.* First we assume, without loss of generality, that at initialisation  $\theta_1 > 0$ , and since  $a$  and  $b$  are not trainable parameters we must have  $a > 0$  and  $b < 0$  to include the ground truth in our hypothesis class.

Let  $J_\epsilon$  be the adversarial training inner maximisation as defined in Equation (39). Then, AT solves the following optimisation problem:

$$\min_{\theta} \mathbb{E}_{(x,y)} [L(\sigma(\text{sgn}(y)J_\epsilon(x,y)), y)] \quad (65)$$

Similarly, let  $\tilde{J}_\epsilon$  be the COAP dual approximation to the inner maximization described in Equation (47). Then, COAP solves the following optimisation problem:

$$\min_{\theta} \mathbb{E}_{(x,y)} [L(\sigma(\text{sgn}(y)\tilde{J}_\epsilon), y)] \quad (66)$$

Since we are only training the signal parameter  $\theta_1$ , after one gradient descent step, we have:

$$\|\bar{\theta}_{2:d}\|_2 = \|\tilde{\theta}_{2:d}\|_2 \quad (67)$$

Further, from Lemma B.1 we know that AT yields smaller robust risk than COAP if the following holds:

$$\|\bar{\theta}_1\|_2 > \|\tilde{\theta}_1\|_2 \implies \mathbf{R}_\epsilon(\tilde{\theta}) > \mathbf{R}_\epsilon(\bar{\theta}) \quad (68)$$

which, after one step of gradient descent, is equivalent to:

$$\mathbb{E}_{(x,y)} [\nabla_{\theta_1} L(\sigma(\text{sgn}(y)J_\epsilon(x,y)), y)] < \mathbb{E}_{(x,y)} [\nabla_{\theta_1} L(\sigma(\text{sgn}(y)\tilde{J}_\epsilon(x,y)), y)]$$

Now recall from Theorems B.1 and B.2 that the gradients of AT and COAP differ only on the set  $\mathcal{I}^*$ . In particular, we have that:

$$(x, y) \notin \mathcal{I}^* \implies \nabla_{\theta_1} L(\sigma(\text{sgn}(y)J_\epsilon(x,y)), y) = \nabla_{\theta_1} L(\sigma(\text{sgn}(y)\tilde{J}_\epsilon(x,y)), y) < 0$$

and

$$(x, y) \in \mathcal{I}^* \implies 0 = \nabla_{\theta_1} L(\sigma(\text{sgn}(y)J_\epsilon(x,y)), y) \neq \nabla_{\theta_1} L(\sigma(\text{sgn}(y)\tilde{J}_\epsilon(x,y)), y)$$

Hence, for our purpose we need to show that:

$$\mathbb{E}_{(x,y)} [\nabla_{\theta_1} L(\sigma(\text{sgn}(y)\tilde{J}_\epsilon(x,y)), y) \mid (x, y) \in \mathcal{I}^*] > 0 \quad (69)$$

Our strategy will be to lower-bound the expectation in Equation (69) with some strictly positive quantity. We define

$$Z = \sum_{i=2}^d \theta_i x_i \quad (70)$$

and plug-in the gradient computed in Theorem B.2:

$$\begin{aligned} & \mathbb{E}_{(x,y)} [\nabla_{\theta_1} L(\sigma(\text{sgn}(y)\tilde{J}_\epsilon(x,y)), y) \mid (x, y) \in \mathcal{I}^*] \\ &= \mathbb{E}_{(x,y)} \left[ \frac{a\sigma(-\tilde{J}_\epsilon(x,y))}{2\epsilon} \left( -\frac{\ell}{\theta_1}(\gamma + \epsilon) + u \frac{\sum_{i=2}^d x_i \theta_i}{\theta_1^2} \right) \mid (x, y) \in \mathcal{I}^* \right] \\ &= \frac{a}{2\theta_1\epsilon} \mathbb{E}_{(x,y)} \left[ \sigma(-\tilde{J}_\epsilon(x,y)) \left( -\ell(\gamma + \epsilon) + u \frac{Z}{\theta_1} \right) \mid (x, y) \in \mathcal{I}^* \right] \\ &= \frac{a}{2\theta_1\epsilon} \mathbb{E}_{(x,y)} \left[ \sigma(-\tilde{J}_\epsilon(x,y)) u \frac{Z}{\theta_1} - \sigma(-\tilde{J}_\epsilon(x,y)) \ell(\gamma + \epsilon) \mid (x, y) \in \mathcal{I}^* \right] \end{aligned} \quad (71)$$

Now, we observe that  $Z$  is always negative on the set  $\mathcal{I}^*$ , since we need to satisfy the constraint  $\ell < 0 < u$ :

$$(x, y) \in \mathcal{I}^* \implies -\theta_1(\gamma + \epsilon) < \sum_{i=2}^d \theta_i x_i < -\theta_1(\gamma - \epsilon) < 0 \quad (72)$$

Further, from Equation (59) we have:

$$(x, y) \in \mathcal{I}^* \implies \sigma\left(-\tilde{J}_\epsilon(x, y)\right) \geq \frac{1}{2} \quad (73)$$

Combining these two observations we can lower-bound the expectation:

$$\begin{aligned} & \mathbb{E}_{(x, y)} \left[ \nabla_{\theta_1} L\left(\text{sgn}(y)\sigma\left(\tilde{J}_\epsilon(x, y)\right), y\right) \mid (x, y) \in \mathcal{I}^* \right] \\ &= \frac{a}{2\theta_1\epsilon} \mathbb{E}_{(x, y)} \left[ \sigma\left(-\tilde{J}_\epsilon(x, y)\right) u \frac{Z}{\theta_1} - \sigma\left(-\tilde{J}_\epsilon(x, y)\right) \ell(\gamma + \epsilon) \mid (x, y) \in \mathcal{I}^* \right] \\ &\geq \frac{a}{2\theta_1\epsilon} \mathbb{E}_{(x, y)} \left[ u \frac{Z}{\theta_1} - \frac{\gamma + \epsilon}{2} \ell \mid (x, y) \in \mathcal{I}^* \right] \end{aligned} \quad (74)$$

Now, we need to show that this lower-bound is strictly positive:

$$\mathbb{E}_{(x, y)} \left[ u \frac{Z}{\theta_1} - \frac{\gamma + \epsilon}{2} \ell \mid (x, y) \in \mathcal{I}^* \right] > 0 \quad (75)$$

Note that, we can further expand this expression:

$$\begin{aligned} & \mathbb{E}_{(x, y)} \left[ u \frac{Z}{\theta_1} - \frac{\gamma + \epsilon}{2} \ell \mid (x, y) \in \mathcal{I}^* \right] \\ &= -(\gamma^2 - \epsilon^2)\theta_1^2 + (\gamma + \epsilon)\theta_1 \mathbb{E}[Z \mid (x, y) \in \mathcal{I}^*] + 2\mathbb{E}[Z^2 \mid (x, y) \in \mathcal{I}^*] \end{aligned} \quad (76)$$

Further,  $Z \mid (x, y) \in \mathcal{I}^*$  is distributed as a truncated normal with:

$$\alpha = -\frac{\theta_1(\gamma + \epsilon)}{\sigma \|\theta_{2:d}\|_2} \quad \text{and} \quad \beta = -\frac{\theta_1(\gamma - \epsilon)}{\sigma \|\theta_{2:d}\|_2} \quad (77)$$

Hence, we can plug in the expectations of the truncated normal distribution to obtain the following:

$$\begin{aligned} & -(\gamma^2 - \epsilon^2)\theta_1^2 + \theta_1(\gamma + \epsilon)\mathbb{E}[Z \mid (x, y) \in \mathcal{I}^*] + 2\mathbb{E}[Z^2 \mid (x, y) \in \mathcal{I}^*] \\ &= -(\gamma^2 - \epsilon^2)\theta_1^2 + 2\sigma^2 \|\theta_{2:d}\|_2^2 + \sigma \|\theta_{2:d}\|_2 \theta_1 \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\Phi(\beta) - \Phi(\alpha)} \\ &\propto -(\gamma^2 - \epsilon^2) + 2\sigma^2 r^2 + \sigma r \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\Phi(\beta) - \Phi(\alpha)} \\ &= -f(r) \end{aligned} \quad (78)$$

where we define  $r = \frac{\|\theta_{2:d}\|_2}{\|\theta_1\|_2}$ . Now, under our assumptions, from Lemma C.3 we have:

$$f(r) < 0, \quad \forall r > \sqrt{\frac{24\gamma^3}{\sigma^2}} \quad (79)$$

which concludes the proof.  $\square$

## C Auxiliary lemmas

### C.1 Upper bound on the exponential function

**Lemma C.1.** *Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be the function defined by  $f(x) = \exp(x)$ . When  $x \leq 0$  and  $n$  is even we have:*

$$f(x) \leq 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!} \quad (80)$$

*Proof.* Let  $g : (-\infty, 0] \rightarrow \mathbb{R}$  be the function defined by

$$g(x) = 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!} - \exp(x) \quad (81)$$

Since  $g(x) \rightarrow \infty$  as  $x \rightarrow -\infty$ ,  $g$  must attain an absolute minimum somewhere on the interval  $(-\infty, 0]$ . Now, differentiating we have:

- If  $f$  has an absolute minimum at 0, then for all  $x$ ,  $f(x) \geq f(0) = 1 - \exp(0) = 0$ , so we are done.
- If  $f$  has an absolute minimum at  $y$  for some  $y < 0$ , then  $f'(y) = 0$ . But differentiating,

$$f'(y) = 1 + y + \frac{y^2}{2!} + \cdots + \frac{y^{n-1}}{(n-1)!} - \exp(y) = f(y) - \frac{y^n}{n!}.$$

Therefore, for any  $x$ ,

$$f(x) \geq f(y) = \frac{y^n}{n!} + f'(y) = \frac{y^n}{n!} > 0,$$

since  $n$  is even.

□

## C.2 Lower bound on the difference of Gaussian CDFs

**Lemma C.2.** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be the function defined by  $f(x, y) = \Phi(y) - \Phi(x)$ . When  $x < y < 0$  we have:

$$\phi(0) \left( y - x + \frac{x^3}{6} \right) \leq \Phi(y) - \Phi(x) \quad (82)$$

where  $\Phi$  and  $\phi$  are respectively the CDF and PDF of the standard Gaussian distribution.

*Proof.* First, we want to prove that  $\frac{2x}{\sqrt{\pi}}$  is a lower bound for the error function  $\operatorname{erf}(x)$  when  $x \leq 0$ . That is, we want to show that  $f(x) \geq 0$  where  $f : (-\infty, 0] \rightarrow \mathbb{R}$  is the function defined by:

$$f(x) = \operatorname{erf}(x) - \frac{2x}{\sqrt{\pi}} \quad (83)$$

Since  $f$  is continuous and  $f(x) \rightarrow \infty$  as  $x \rightarrow -\infty$ ,  $f$  must attain an absolute minimum on the interval  $(-\infty, 0]$ . Now, differentiating we have:

$$f'(x) = \frac{2}{\sqrt{\pi}} \exp(-x^2) - \frac{2}{\sqrt{\pi}} \quad (84)$$

hence  $f$  attains an absolute minimum at 0 and we have  $f(x) \geq f(0) = 0$ .

Next, we show that  $\frac{2}{\sqrt{\pi}}(x - x^3/3)$  is an upper bound for  $\operatorname{erf}(x)$  when  $x \leq 0$ . Let  $g : (-\infty, 0] \rightarrow \mathbb{R}$  the function defined by:

$$g(x) = \frac{2}{\sqrt{\pi}}(x - x^3/3) - \operatorname{erf}(x) \quad (85)$$

Similarly, since  $g$  is continuous and  $g(x) \rightarrow \infty$  as  $x \rightarrow -\infty$ ,  $g$  must attain an absolute minimum on the interval  $(-\infty, 0]$ . Now, differentiating we have:

$$g'(x) = \frac{2}{\sqrt{\pi}}(1 - x^2 - \exp(-x^2)) \quad (86)$$

hence  $g$  attains an absolute minimum at 0 and we have  $g(x) \geq g(0) = 0$ . Now, since  $a < b < 0$  we can use the erf bounds derived above:

$$\Phi(b) - \Phi(a) = \frac{1}{2} \left( \operatorname{erf}(b/\sqrt{2}) - \operatorname{erf}(a/\sqrt{2}) \right) \quad (87)$$

$$\geq \frac{1}{\sqrt{\pi}} \left( \frac{b}{\sqrt{2}} - \frac{a}{\sqrt{2}} + \frac{a^3}{6\sqrt{2}} \right) \quad (88)$$

$$= \phi(0) \left( b - a + \frac{a^3}{6} \right) \quad (89)$$

which concludes the proof. □

### C.3 Upper bound on the ratio of Gaussian PDFs and CDFs

**Lemma C.3.** Suppose  $f : \mathbb{R} \rightarrow \mathbb{R}$  is defined as follows:

$$f(r) = \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma r \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\Phi(\beta) - \Phi(\alpha)} \quad (90)$$

where  $\alpha := -\frac{\gamma+\epsilon}{r\sigma}$ ,  $\beta := -\frac{\gamma-\epsilon}{r\sigma}$ ,  $\Phi$  and  $\phi$  are respectively the standard Gaussian CDF and PDF.

Assume that:

$$\frac{5 + 2\sqrt{3}}{13}\gamma < \epsilon < \gamma \quad (91)$$

Then, we have:

$$f(r) < 0, \quad \forall r > \sqrt{\frac{24\gamma^3}{\sigma^2}} \quad (92)$$

*Proof.* We begin by providing a lower bound on the difference of gaussian cdfs. Applying Lemma C.2 with  $x = \alpha$  and  $y = \beta$  we have:

$$\Phi(\beta) - \Phi(\alpha) \geq \left( \frac{2\epsilon}{r\sigma} - \frac{(\gamma + \epsilon)^3}{6\sigma^3 r^3} \right) \phi(0), \quad \alpha < \beta < 0 \quad (93)$$

Next, we can upper-bound  $f$ :

$$f(r) \leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma r \frac{(\gamma - 3\epsilon)\phi(\beta) - (\gamma + \epsilon)\phi(\alpha)}{\left( \frac{2\epsilon}{r\sigma} - \frac{(\gamma + \epsilon)^3}{6\sigma^3 r^3} \right) \phi(0)} \quad (94)$$

$$\leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon)\phi(0) - (\gamma + \epsilon)\phi(\alpha)}{\left( 2\epsilon - \frac{(\gamma + \epsilon)^3}{6r^2\sigma^2} \right) \phi(0)} \quad (95)$$

$$= \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon) - (\gamma + \epsilon)\exp(-\alpha^2/2)}{2\epsilon - \frac{(\gamma + \epsilon)^3}{6\sigma^2 r^2}} \quad (96)$$

Now, we use the upper-bound for the exponential function from Lemma C.1 with  $n = 2$ :

$$\exp(x) \leq 1 + x - x^2/2, \quad \forall x \leq 0 \quad (97)$$

and substituting it back into our upper-bound for  $f$  we get:

$$f(r) \leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon) - (\gamma + \epsilon)\left(1 - \frac{(\gamma + \epsilon)^2}{2r^2\sigma^2} + \frac{(\gamma + \epsilon)^4}{8r^4\sigma^4}\right)}{2\epsilon - \frac{(\gamma + \epsilon)^3}{6r^2\sigma^2}} \quad (98)$$

which can be further simplified:

$$f(r) \leq \gamma^2 - \epsilon^2 - 2\sigma^2 r^2 - \sigma^2 r^2 \frac{(\gamma - 3\epsilon) - (\gamma + \epsilon)\left(1 - \frac{(\gamma + \epsilon)^2}{2r^2\sigma^2} + \frac{(\gamma + \epsilon)^4}{8r^4\sigma^4}\right)}{2\epsilon - \frac{(\gamma + \epsilon)^3}{6r^2\sigma^2}} \quad (99)$$

$$= \frac{(\gamma - 7\epsilon)(\gamma + \epsilon)^4 + 4r^2\sigma^2(\gamma + \epsilon)(\gamma^2 - 10\gamma\epsilon + 13\epsilon^2)}{4(\gamma + \epsilon)^3 - 48r^2\sigma^2\epsilon} \quad (100)$$

$$= u(r) \quad (101)$$

and we have that for  $\epsilon > \frac{5+2\sqrt{3}}{13}\gamma$  and  $r > \sqrt{\max\left(\frac{(7\epsilon-\gamma)(\gamma+\epsilon)^4}{4\sigma^2(\gamma^2-10\gamma\epsilon+13\epsilon^2)}, \frac{(\gamma+\epsilon)^3}{12\sigma^2\epsilon}\right)}$  the upper bound is negative, i.e.  $u(r) < 0$ . Finally, for the sake of clarity we can further simplify the condition on  $r$ :

$$r > \sqrt{\frac{24\gamma^3}{\sigma^2}} > \sqrt{\max\left(\frac{(7\epsilon-\gamma)(\gamma+\epsilon)^4}{4\sigma^2(\gamma^2-10\gamma\epsilon+13\epsilon^2)}, \frac{(\gamma+\epsilon)^3}{12\sigma^2\epsilon}\right)} \quad (102)$$

which concludes the proof.  $\square$

## D Experimental details

### D.1 Synthetic experiments with signal-directed adversaries

Below we provide detailed experimental details to reproduce Figure 5.

**Data generation** For the linearly separable distribution we set  $d = 1000$ ,  $n_{\text{test}} = 10^5$ ,  $\gamma = 6$ .

**Model and hyper-parameters** For all the experiments, we use the one hidden layer architecture defined in Equation (8) with 100 neurons. We use PyTorch SGD optimiser and train all networks for 100 epochs. We sweep over the learning rate  $\eta \in \{0.1, 0.01, 0.001\}$  and for each perturbation budget, we choose the one that interpolates the training set and minimises robust error on the test set.

**Robust evaluation** We perform all the attacks to evaluate robust risk at test-time using exact line search; this is computationally tractable since the attacks are directed along one dimension.

**Training paradigms** For standard training (ST), we train the network to minimise the cross-entropy loss. For adversarial training (AT) [20, 13], we train the network to minimise the robust binary cross-entropy loss. At each epoch, we compute an exact adversarial example using line search and update the weights using a gradient with respect to this example. For convex outer adversarial polytope (COAP) [31, 32], at each epoch, we compute upper and lower bounds  $u$  and  $\ell$  as described in Theorem A.1. We then train the network to minimise the upper bound on robust error from Theorem A.2.

**Standard training.** We train the network to minimise the cross-entropy loss.

**Adversarial training** [20, 13]. We train the network to minimise the robust binary cross-entropy loss. At each epoch, we compute an exact adversarial example using line search and update the weights using a gradient with respect to this example.

**Certified training** [31, 32]. At each epoch, we compute upper and lower bounds  $u$  and  $\ell$  as described in Proposition A.1. We then train the network to minimize the upper-bound on robust error derived in Theorem A.2.

### D.2 Synthetic experiments with $\ell_2$ adversaries

Below we provide complete experimental details to reproduce Figures 3 and 4.

**Data generation** For the spheres dataset, we generate a random  $x \in \mathbb{R}^d$  where  $\|x\|_2$  is either  $R_1$  or  $R_{-1}$ , with equal probability assigned to each norm. We associate with each  $x$  a label  $y$  such that  $y = -1$  if  $\|x\|_2 = R_{-1}$  and  $y = 1$  if  $\|x\|_2 = R_1$ . We can sample uniformly from this distribution by sampling  $z \sim \mathcal{N}(0, I_d)$  and then setting  $x = \frac{z}{\|z\|_2} R_{-1}$  or  $x = \frac{z}{\|z\|_2} R_1$ . For the linearly separable distribution we set  $d = 1000$ ,  $n = 50$ ,  $n_{\text{test}} = 10^5$ ,  $\gamma = 6$ . For the concentric spheres distribution we set  $d = 100$ ,  $n = 50$ ,  $n_{\text{test}} = 10^5$ ,  $\gamma_{\min} = 1$  and  $\gamma_{\max} = 12$ .

**Model and hyper-parameters** For all the experiments, we use a MLP architecture with  $W = 100$  neurons in each hidden layer and ReLU( $\cdot$ ) activation functions. We use PyTorch SGD optimiser with a momentum of 0.95 and train the network for 150 epochs. We sweep over the learning rate  $\eta \in \{0.1, 0.01, 0.001\}$  and for each perturbation budget, we choose the one that minimises robust error on the test set and interpolates the training set.

**Robust evaluation** We consider  $\ell_2$ -ball perturbations. We evaluate robust error at test-time using Auto-PGD [4] with 100 iterations and 5 random restarts. We use both the cross-entropy and difference of logits loss to prevent gradient masking. We use the implementation provided in AutoAttack [4] with minor adjustments to allow for non-image inputs.

**Training paradigms** For standard training (ST), we train the network to minimise the cross-entropy loss. For adversarial training (AT) [20, 13] we train the network to minimise the robust cross-entropy loss. At each epoch, we search for adversarial examples using Auto-PGD [4] with a budget of 10 steps

and 1 random restart. Then, we update the weights using a gradient with respect to the adversarial examples. For convex outer adversarial polytope (COAP) [31, 32]. We train the network to minimise the upper-bound on the robust error. Our implementation is based on the code released by the authors.

### D.3 Image experiments

Below we provide complete experimental details to reproduce Figure 2.

**Model architectures** For MNIST, we train the CNN architecture with four convolutional layer and two fully connected layers of 512 units introduced in Wong et al. [32]. We report the architectural details in Table 1. For CIFAR-10, we train the residual network (ResNet) with the same structure used in Wong et al. [32]; we use 1 residual block with 16, 16, 32, and 64 filters. For Tiny ImageNet, we train a WideResNet. Following Xu et al. [34] we use 3 wide basic blocks with a widen factor of 10.

CNN
CONV 32 $3 \times 3 + 1$
CONV 32 $4 \times 4 + 2$
CONV 64 $3 \times 3 + 1$
CONV 64 $4 \times 4 + 2$
FC 512
FC 512

Table 1: MNIST model architecture. All layers are followed by ReLU ( $\cdot$ ) activations. The last fully connected layer is omitted. "CONV  $k w \times h + s$ " corresponds to a 2D convolutional layer with  $k$  filters of size  $w \times h$  using a stride of  $s$  in both dimensions. "FC  $n$ " corresponds to a fully connected layer with  $n$  outputs.

**Dataset preprocessing** For MNIST, we use full  $28 \times 28$  images without any augmentations and normalisation. For CIFAR-10, we use random horizontal flips and random crops as data augmentation, and normalise images according to per-channel statistics. For Tiny ImageNet, we use random crops of  $56 \times 56$  and random flips during training. During testing, we use a central  $56 \times 56$  crop. We also normalise images according to per-channel statistics.

**Robust evaluation** We consider  $\ell_2$ -ball perturbations. We evaluate the robust error using the most expensive version of AutoAttack (AA+) [4]. Specifically, we include the following attacks: untargeted APGD-CE (5 restarts), untargeted APGD-DLR (5 restarts), untargeted APGD-DLR (5 restarts), Square Attack (5000 queries), targeted APGD-DLR (9 target classes) and targeted FAB (9 target classes).

**AT training details** For MNIST, we train 100 epochs using Adam optimiser [17] with a learning rate of 0.001, momentum of 0.9 and a batch size of 128; we reduce the learning rate by a factor 0.1 at epochs 40 and 80. For CIFAR-10 with ResNet, we train 150 epochs using SGD with a learning rate of 0.05 and a batch size of 128; we reduce the learning rate by a factor 0.1 at epochs 80 and 120. For Tiny Imagenet and CIFAR-10 with Wide-Resnet we train 200 epochs using SGD with a learning rate of 0.1 and a batch size of 512; we reduce the learning rate by a factor 0.1 at epochs 100 and 150. For the inner optimisation of all models and datasets, adversarial examples are generated with 10 iterations of Auto-PGD [4].

**COAP training details** We follow the settings proposed by the authors and report them here. For MNIST, we use the Adam optimiser [17] with a learning rate of 0.001 and a batch size of 50. We schedule  $\epsilon$  starting from 0.01 to the desired value over the first 20 epochs, after which we decay the learning rate by a factor of 0.5 every 10 epochs for a total of 60 epochs. For CIFAR-10, we use the SGD optimiser with a learning rate of 0.05 and a batch size of 50. We schedule  $\epsilon$  starting from 0.001 to the desired value over the first 20 epochs, after which we decay the learning rate by a factor of 0.5 every 10 epochs for a total of 60 epochs. For all datasets and models, we use random projection of 50 dimensions. For all experiments, we use the implementation provided in Wong et al. [32].

**CROWN-IBP training details** We follow the settings proposed by the authors and report them here. For MNIST, we train 200 epochs with a batch size of 256. We use Adam optimiser [17] and set the learning rate to  $5 \times 10^{-4}$ . We warm up with 10 epochs of regular training, and gradually ramp up  $\epsilon_{\text{train}}$  from 0 to  $\epsilon$  in 50 epochs. We reduce the learning rate by a factor 0.1 at epoch 130 and 190. For CIFAR-10, we train 2000 epochs with a batch size of 256, and a learning rate of  $5 \times 10^{-4}$ . We warm up for 100 epochs, and ramp-up  $\epsilon$  for 800 epochs. Learning rate is reduced by a factor 0.1 at epoch 1400 and 1700. For Tiny ImageNet, we train 600 epochs with batch size 128. The first 100 epochs are clean training, then we gradually increase  $\epsilon_{\text{train}}$  with a schedule length of 400. For all datasets, an hyper-parameter  $\beta$  to balance LiRPA bounds and IBP bounds for the output layer is gradually decreased from 1 to 0 (1 for only using LiRPA bounds and 0 for only using IBP bounds), with the same schedule of  $\epsilon$ . For all experiments, we use the implementation provided in the auto LiRPA library [34].