045

046

053

054

058

# The Price of Robustness: Stable Classifiers Need Overparameterization

Jonas von Berg\*1,2, Adalbert Fono<sup>1,2</sup>, Massimiliano Datres<sup>1,2</sup>, Sohir Maskey<sup>1,2,3</sup>, and Gitta  $Kutyniok^{1,2,4,5}$ 

{berg, fono, datres, kutyniok}@math.lmu.de

sohir.maskey@aleph-alpha-research.com

# Abstract

002

003

006

007

008

009

010

011

014

015

016

017

018

020

021

022

023

024

025

028

029

030

031

032

035

036

037

038

The link between overparameterization, robustness, and generalization in discontinuous classifiers remains unclear. We establish generalization bounds that tighten with class stability – the expected distance to the decision boundary – yielding a law of robustness for classification that extends prior smoothness based settings. As a consequence, any interpolating model with  $p \approx n$  parameters is necessarily unstable, implying that robust generalization requires overparameterization. For infinite function classes, we obtain analogous results through a stronger robustness measure, the normalized costability, defined via output margins. Empirical results support our theory: stability grows with model size and aligns closely with test performance.

#### Introduction 1

The generalization behavior of overparameterized neural networks presents fundamental challenges to classical statistical learning theory. Traditional complexity measures, such as parameter counts or spectral norms of weights, form the basis of many generalization bounds, including those derived from VC dimension theory [1] and Rademacher complexity [2]. However, these approaches do not adequately explain several empirical phenomena, e.g., double descent [3] and benign overfitting [4], where test performance improves beyond the interpolation threshold. Empirical studies further show that norm-based metrics often correlate poorly with generalization [5], while the margin – the distance to the decision boundary – emerges as a reliable predictor [6–8]. This suggests that generalization is governed not by microscopic weight norms but by macroscopic simplicity, the stability of predictions under perturbations. The law of robustness of Bubeck and Sellke [9] establishes a formal link between robustness, generalization, and overparameterization: smoothness and overparameterization need to balance in order to ensure good generalization while overfitting. However, its reliance on smoothness assumptions excludes discontinuous classifiers. We address this limitation by introducing class stability and normalized co-stability – geometric, macroscopic measures of functional simplicity that extend robustness laws to classification.

### 2 Setup

We study binary classification on  $(\mathcal{X} \times \{-1, 1\}, \mu)$ , 048 where  $\mathcal{X} \subset \mathbb{R}^d$  is bounded and  $\mathcal{F} \subset \{f : \mathcal{X} \to \{-1, 1\}\}$ a hypothesis class. Given n i.i.d. samples  $(x_i, y_i) \sim \mu$ , the goal is to find  $f \in \mathcal{F}$  minimizing a bounded loss  $\ell$ . We focus on the binary case; multi-class extensions follow by one-vs-all reduction (A.4). A canonical loss is  $\ell_{0-1}(y, y') = \mathbb{1}_{y \neq y'}$ .

Class stability. Following Liu and Hansen [10], we measure robustness by the expected distance to the decision boundary. For  $f: \mathcal{X} \to \{-1, 1\}$ , define the signed distance function

$$d_f(x) = \begin{cases} d(x, f^{-1}(\{-1\})), & f(x) = 1, \\ -d(x, f^{-1}(\{1\})), & f(x) = -1, \end{cases}$$
(1) 059

where  $d(x, A) = \inf_{y \in A} ||x - y||_2$ . The (unsigned) 060 margin and the class stability are 061

$$h_f(x) = |d_f(x)|, \qquad S(f) = \mathbb{E}[h_f].$$
 (2) 062

Here S(f) quantifies the average distance of samples to the decision boundary – a notion of macroscopic robustness. To extend our results to infinite, parameterized function classes, continuity in the parameterization is required. For this, we introduce a stronger, codomain-based notion of stability. 068

**Co-stability.** Any classifier can be written as f = $\operatorname{sgn} \circ g$ , where g is Lipschitz continuous with constant L(q) (see Lemma 2). This representation allows us

<sup>&</sup>lt;sup>1</sup>Ludwig-Maximilians-Universität München

<sup>&</sup>lt;sup>2</sup> Munich Center for Machine Learning (MCML)

<sup>&</sup>lt;sup>3</sup>Aleph Alpha Research

<sup>&</sup>lt;sup>4</sup>University of Tromsø

<sup>&</sup>lt;sup>5</sup>DLR-German Aerospace Center

<sup>\*</sup>Corresponding Author.

075

077

080

081

082

083

084

086

087

088

089

090

091

092

093

095

096

097

098

101

102

107

108

109

110

112

113

114

119

120

121

127

128

131

132

138

147

148

to define the (Lipschitz-normalized) co-margin and co-stability as 073

$$\bar{h}_g^*(x) = \frac{|g(x)|}{L(g)}, \qquad \bar{S}^*(g) = \mathbb{E}[\bar{h}_g^*(x)].$$
 (3)

For the canonical choice  $g = d_f$ , they coincide with S(f). Since perturbing x by r changes g(x) by at most L(g)r, label flips require  $r \ge |g(x)|/L(g)$ , implying that in general we have the inequality

$$S(f) \ge \bar{S}^*(g). \tag{4}$$

Normalized co-stability thus lower-bounds class stability and provides another scale-invariant robustness measure. Crucially, it guarantees that the score function g remains, on average, at a nontrivial distance from the decision discontinuity, enabling extensions to infinite function classes.

**Isoperimetry.** To control how stable functions fit random labels, we assume  $\mu$  satisfies concentration for Lipschitz functions:

$$\mathbb{P}(|f(x) - \mathbb{E}[f]| \ge t) \le 2e^{-dt^2/(2cL^2)}$$
 (5)

for all bounded L-Lipschitz  $f: \mathcal{X} \to \mathbb{R}$  and  $t \geq 0$ . This c-isoperimetry holds for Gaussian measures and for uniform measures on compact manifolds with positive curvature [9, 11]. Under the manifold hypothesis, d represents intrinsic dimension.

## 3 Main Results: A Law of Robustness for Classification

We establish a law of robustness for classification, linking generalization to margin-based stability in discontinuous classifiers. In contrast to Lipschitzbased analyses [9], our bounds apply directly to discrete decision functions through the notions of class stability S(f) and normalized co-stability  $\bar{S}^*(g)$ .

Finite class Rademacher bound. If the input 103 distribution is c-isoperimetric and  $\mathcal{F}$  is finite with  $\min_{f \in \mathcal{F}} S(f) \geq S > 0$  and  $\log |\mathcal{F}| \geq n$ , then the 105 Rademacher complexity satisfies 106

$$\mathcal{R}_{n,\mu}(\mathcal{F}) \lesssim \max \left\{ \frac{1}{\sqrt{n}}, \frac{1}{S} \sqrt{\frac{c \log |\mathcal{F}|}{n d}} \right\}.$$
 (6)

(Precise statement and proof: Theorem 4.)

Infinite class Rademacher bound. For parameterized classes  $f = \operatorname{sgn} \circ g_{\omega}$  with bounded parameter set  $\mathcal{W} \subset \mathbb{R}^p$ , where  $g_{\omega}$  is Lipschitz in  $\mathcal{X}$   $(L_{\mathcal{X}}(g) \leq L)$ and Lipschitz in  $\omega$ , and  $\bar{S}^*(g_\omega) \geq S^* > 0$ , one ob-

$$\mathcal{R}_{n,\mu}(\mathcal{F}) \lesssim \max \left\{ \frac{1}{\sqrt{n}}, \frac{L}{S^*} \sqrt{\frac{cp}{nd}} \right\}.$$
 (7)

(Precise statement and proof: Theorem 5.)

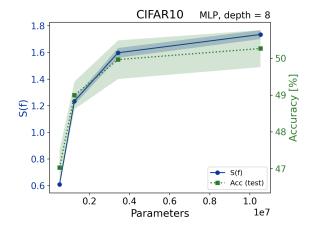


Figure 1. Class stability for MLPs trained on CIFAR-

Law of Robustness. Combining the above results with the standard generalization bound in terms of the Rademacher complexity [2] yields the following informal statement. If  $\sigma^2 := \min_{f \in \mathcal{F}} R_{0-1}(f) >$  $\varepsilon > 0$  and a classifier satisfies  $\hat{R}_{0-1}(f) \leq \sigma^2 - \varepsilon$  for sufficiently large n, then with high probability

$$\frac{S^*(g)}{L(g)} \lesssim \frac{1}{\varepsilon} \sqrt{\frac{c\,p}{n\,d}}.\tag{8}$$

Hence, simultaneously achieving low training error and high (co-)stability requires overparameterization on the order of  $p \approx nd$ . An analogous relation holds for finite function classes in terms of S(f). (Precise finite and infinite formulations, together with proofs, are given in the appendix A.3.)

**Experiments.** We trained 4- and 8-layer MLPs of varying width on MNIST and CIFAR-10, estimating S(f) via minimal  $\ell_2$  adversarial radii and  $\bar{S}^*$  via efficient lipschitz estimation (using the ECLIPSE method [12]). Both measures increase with width and correlate strongly with test accuracy. These trends, support our theory that (co-)stability, grows with overparameterization. Experimental details and further plots are provided in the appendix A.5.

#### 4 Conclusion

Our results show that good generalization in over- 139 parameterized regimes hinges on sufficient stability. 140 The inverse dependence on S or  $\bar{S}^*/L$  in our bounds indicates that stability reduces effective complexity, mitigating overfitting. In high dimensions, overparameterization becomes necessary for robust generalization: limited capacity forces a trade-off with (co-)stability, leading to large Lipschitz constants or low prediction confidence. This aligns with observations that large neural networks, including LLMs, generalize well despite overparameterization.

210

213

214

216

217

218

220

224

239

240

244

245

247

250

251

## References

- S. R. Sain. The nature of statistical learning 151 theory. 1996. 152
- P. L. Bartlett and S. Mendelson. "Rademacher 153 and gaussian complexities: Risk bounds and structural results". In: Journal of Machine 155 Learning Research 3.Nov (2002), pp. 463–482. 156
- 157 M. Belkin, D. Hsu, S. Ma, and S. Mandal. "Reconciling modern machine-learning practice 158 and the classical bias-variance trade-off". In: 159 Proceedings of the National Academy of Sci-160 ences 116.32 (2019), pp. 15849–15854. 161
- P. L. Bartlett, P. M. Long, G. Lugosi, and A. 162 163 Tsigler. "Benign overfitting in linear regression". In: Proceedings of the National Academy 164 of Sciences 117.48 (2020), pp. 30063–30070. 165
- Y. Jiang, B. Nevshabur, H. Mobahi, D. Kr-166 ishnan, and S. Bengio. Fantastic Generaliza-167 tion Measures and Where to Find Them. 2019. arXiv: 1912.02178 [cs.LG]. URL: https:// 169 arxiv.org/abs/1912.02178. 170
- B. Zhang, D. Jiang, D. He, and L. Wang. Re-171 thinking Lipschitz Neural Networks and Certi-172 fied Robustness: A Boolean Function Perspec-173 tive. 2022. arXiv: 2210.01787 [cs.LG]. URL: 174 https://arxiv.org/abs/2210.01787. 175
- N. Ghosh and M. Belkin. A Universal Trade-176 off Between the Model Size, Test Loss, and 177 Training Loss of Linear Predictors. 2023. 178 arXiv: 2207.11621 [stat.ML]. URL: https: //arxiv.org/abs/2207.11621.
- J. A. Soloff, R. F. Barber, and R. Willett. 181 Building a stable classifier with the inflated argmax. 2025. arXiv: 2405.14064 [stat.ML]. 183 URL: https://arxiv.org/abs/2405.14064. 184
- S. Bubeck and M. Sellke. "A universal law 185 of robustness via isoperimetry". In: Advances 186 in Neural Information Processing Systems 34 187 (2021), pp. 28811-28822. 188
- Z. N. D. Liu and A. C. Hansen. Do stable neu-[10]189 ral networks exist for classification problems? 190 - A new view on stability in AI. 2024. arXiv: 191 2401.07874 [cs.LG]. URL: https://arxiv. 192 org/abs/2401.07874. 193
- R. Vershynin. High-Dimensional Probability: 194 An Introduction with Applications in Data Sci-195 ence. Cambridge Series in Statistical and Prob-196 abilistic Mathematics. Cambridge University Press, 2018. 198

- Y. Xu and S. Sivaranjani. "ECLipsE: Ef- 199 ficient Compositional Lipschitz Constant Estimation for Deep Neural Networks". In: Advances in Neural Information Pro- 202 cessing Systems. Ed. by A. Globerson, L. 203 Mackey, D. Belgrave, A. Fan, U. Paquet, 204 J. Tomczak, and C. Zhang. Vol. 37. Curran 205 Associates, Inc., 2024, pp. 10414–10441. 206 URL: https://proceedings.neurips.207 cc / paper \_ files / paper / 2024 / file / 1419d8554191a65ea4f2d8e1057973e4 Paper-Conference.pdf.
- [13] Rigollet J.-C. Ρ. and Hütter. High- 211 Dimensional 2023. arXiv: 212 Statistics. [math.ST]. 2310 . 19244 URL: https: //arxiv.org/abs/2310.19244.
- [14]D. P. Kingma and J. Ba. "Adam: A Method for Stochastic Optimization". In: 3rd International Conference on Learning Representations (ICLR). 2015. URL: https://arxiv.org/abs/ 1412.6980.
- [15]K. Jordan, Y. Jin, V. Boza, J. You, F. Cesista, L. Newhouse, and J. Bernstein. Muon: An optimizer for hidden layers in neural networks. 2024. URL: https://kellerjordan.github. 223 io/posts/muon/.
- [16] P. Foret, A. Kleiner, H. Mobahi, and B. 225 Neyshabur. "Sharpness-aware Minimization 226 for Efficiently Improving Generalization". In: 227 International Conference on Learning Representations. 2021. URL: https://openreview. 229 net/forum?id=6Tm1mposlrM.
- J. Kwon, J. Kim, H. Park, and I. K. Choi. 231 [17]"ASAM: Adaptive Sharpness-Aware Minimization for Scale-Invariant Learning of Deep Neural Networks". In: Proceedings of the 38th International Conference on Machine Learning. Ed. by M. Meila and T. Zhang. Vol. 139. 236 Proceedings of Machine Learning Research. 237 PMLR, 18–24 Jul 2021, pp. 5905–5914. URL: https://proceedings.mlr.press/v139/ kwon21b.html.
- I. J. Goodfellow, J. Shlens, and C. Szegedy. 241 "Explaining and Harnessing Adversarial Examples". In: International Conference on Learning Representations (ICLR) (2015). URL: https://arxiv.org/abs/1412.6572.
- S.-M. Moosavi-Dezfooli, A. Fawzi, and P. 246 Frossard. "DeepFool: A Simple and Accurate Method to Fool Deep Neural Networks". In: 248 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 2016, pp. 2574-2582. DOI: 10.1109/CVPR.2016.282.

308

321

322

331

341

- 252 [20] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. "Towards Deep Learning Models Resistant to Adversarial Attacks". In: *International Conference on Learning Representations (ICLR)*. 2018. URL: https://arxiv.org/abs/1706.06083.
- [21] J. Rauber, W. Brendel, and M. Bethge. "Foolbox: A Python toolbox to benchmark the robustness of machine learning models". In:
  arXiv preprint arXiv:1707.04131 (2017). URL: https://arxiv.org/abs/1707.04131.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. [22]263 Bradbury, G. Chanan, T. Killeen, Z. Lin, 264 N. Gimelshein, L. Antiga, A. Desmaison, A. 265 266 Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, 267 J. Bai, and S. Chintala. "PyTorch: An Im-268 perative Style, High-Performance Deep Learn-269 ing Library". In: Advances in Neural Infor-270 mation Processing Systems 32. Curran Asso-271 ciates, Inc., 2019, pp. 8024-8035. URL: http:// 272 papers.neurips.cc/paper/9015-pytorch-273 an-imperative-style-high-performance-274 deep-learning-library.pdf. 275
- [23] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness May Be at
  Odds with Accuracy. 2019. arXiv: 1805.12152
  [stat.ML]. URL: https://arxiv.org/abs/
  1805.12152.
- 281 [24] H. Zhang, Y. Yu, J. Jiao, E. P. Xing, L. E. Ghaoui, and M. I. Jordan. Theoretically Principled Trade-off between Robustness and Accuracy. 2019. arXiv: 1901.08573 [cs.LG]. URL: https://arxiv.org/abs/1901.08573.
- 286 [25] L. Béthune, T. Boissin, M. Serrurier, F. Ma287 malet, C. Friedrich, and A. González-Sanz.
  288 Pay attention to your loss: understanding
  289 misconceptions about 1-Lipschitz neural net290 works. 2022. arXiv: 2104.05097 [cs.LG]. URL:
  291 https://arxiv.org/abs/2104.05097.

# A Appendix

292

## 293 A.1 The Signed Distance Function

We collect the main properties of the signed distance function

$$d_f(x) := \begin{cases} d(x, f^{-1}(\{-1\})), & \text{if } f(x) = 1, \\ -d(x, f^{-1}(\{1\})), & \text{if } f(x) = -1, \end{cases}$$
 (9)

where  $d(x, A) := \inf_{y \in A} ||x - y||_2$ .

**Lemma 1.** Let  $\mathcal{X} \subset \mathbb{R}^d$  be bounded and pathconnected, and let  $f: \mathcal{X} \to \{-1,1\}$ . Then the signed distance function  $d_f$  is 1-Lipschitz. This is a classical fact, a special case of the Eikonal equation. For completeness, we include a direct proof inspired by Liu and Hansen [10, Prop. 7.5].

*Proof.* Case 1: f(x) = f(y). Assume w.l.o.g. 304 f(x) = f(y) = 1. Let  $(z_n)_n$  be a sequence in 305  $f^{-1}(\{-1\})$  with  $|d(y, z_n) - d_f(y)| \le \frac{1}{n}$ . Then 306

$$d_f(x) = d(x, f^{-1}(\{-1\})) \tag{10}$$

$$\leq d(x, z_n) \tag{11}$$

$$\leq \|x - y\|_2 + d(y, z_n)$$
 (12) 309

$$\leq ||x - y||_2 + d_f(y) + \frac{1}{n}.$$
 (13) 310

Letting  $n \to \infty$  and exploiting symmetry yields 311  $|d_f(x) - d_f(y)| \le ||x - y||_2$ . 312

Case 2:  $f(x) \neq f(y)$ . Assume w.l.o.g. f(x) = 1, 313 f(y) = -1. Consider the line segment  $L = \{(1 - 314 t)x + ty : t \in [0, 1]\} \subset \mathcal{X}$  and define 315

$$w_1 = (1 - t_1)x + t_1y, (14)$$

$$t_1 := \inf\{t : f((1-t)x + ty) = -1\},$$
 (15) 317

$$w_2 = (1 - t_2)x + t_2y, (16) 31$$

$$t_2 := \sup\{t : f((1-t)x + ty) = 1\}.$$
 (17) 319

Path-connectedness ensures  $t_1 \leq t_2$ , since otherwise the midpoint between  $w_1$  and  $w_2$  would be labeled both 1 and -1, a contradiction.

 $|d_f(x) - d_f(y)| = d(x, f^{-1}(\{-1\})) + d(y, f^{-1}(\{1\}))$ (18)

$$\leq \|x - w_1\|_2 + \|y - w_2\|_2 \tag{19}$$

$$\leq \|x - y\|_2.$$
 (20) 326

**Lemma 2.** Let  $\mathcal{X} \subset \mathbb{R}^d$  and  $f: \mathcal{X} \to \{-1, 1\}$  with 328  $f^{-1}(\{1\})$  closed. Then f can be represented as 329

$$f(x) = \operatorname{sgn}(d_f(x)), \tag{21}$$

where we adopt the convention sgn(0) = 1.

*Proof.* If  $d_f(x) \neq 0$ , the claim follows directly from the definition of  $d_f$ . If  $d_f(x) = 0$ , then  $x \in f^{-1}(\{1\})$  333 by closedness, so  $f(x) = 1 = \operatorname{sgn}(0)$ .

**Remark 3.** Lemma 2 justifies the representation  $f = \operatorname{sgn} \circ d_f$  used in the proof of Theorem 4. This link between classifiers and their signed distance functions is what allows stability arguments to be combined with smoothness-based tools.

# A.2 Proofs of the Rademacher Bounds

We now provide proofs for the Rademacher bounds of for finite and infinite function classes.

360

361

362

364

365

366

367

370

371

372

373

376

383

381

385

387

388

403

408

Finite Rademacher Bound. We begin by restating the assumptions.

- 346 (H1)  $(\mathcal{X}, \mu)$  is a probability space with bounded sample space  $\mathcal{X}$  and c-isoperimetric measure  $\mu$ ;
- 348 (H2) the considered hypothesis class  $\mathcal{F}$  of classifiers  $f: \mathcal{X} \to \{-1, 1\}$  is finite, that is  $|\mathcal{F}| < \infty$ .

Theorem 4 (Finite Rademacher Bound). Suppose Assumptions (H1) and (H2) hold, and that  $\min_{f \in \mathcal{F}} S(f) > S > 0$  with  $\log |\mathcal{F}| \geq n$ . Let us furthermore assume that  $f^{-1}(\{1\})$  is closed and  $\mathcal{X}$  path connected, then the empirical Rademacher complexity satisfies

356 
$$\mathcal{R}_{n,\mu}(\mathcal{F}) \leq K_2 \max \left\{ \frac{1}{\sqrt{n}}, \frac{\sqrt{c}}{S} \sqrt{\frac{\log |\mathcal{F}|}{nd}}, \right.$$
357 
$$2 \exp\left(-\frac{dS^2}{8c}\right) \right\}. \quad (22)$$

for an absolute constant  $K_2 > 0$ .

Proof. By Lemma 2, every f admits the representation  $f = \operatorname{sgn} \circ d_f$ . This allows us to follow the infinite-class analysis (see the proof of Theorem 5) without requiring the  $\varepsilon$ -net construction in Equation 30. By Lemma 1, the signed distance function  $d_f$  is 1-Lipschitz, i.e.,  $L(d_f) = 1$  under the stated conditions. Moreover, recalling the definition of co-stability, we obtain

$$S^*(d_f) = \mathbb{E}[|d_f|] = \mathbb{E}[h_f] = S(f).$$
 (23)

Plugging this into the general bound in Equation 25 gives the result.

Infinite Rademacher bound We extend the finite-class result to infinite function classes via a covering-number argument, for which the Lipschitz continuity of the parameterization plays a crucial role. To this end, we introduce a new regularity assumption that replaces the finiteness condition (H2).

377 (H3) The hypothesis class  $\mathcal{F}$  is of the form  $\mathcal{F} =$  sgn  $\circ \mathcal{G}$ , where  $\mathcal{G} = \{g_w : \mathcal{X} \to [-1,1] : w \in \mathcal{W}\}$  is a parameterized class of Lipschitz continuous functions. The parameter space  $\mathcal{W} \subset \mathbb{R}^p$  is bounded with diam $(\mathcal{W}) \leq W$ , and the parameterization is Lipschitz continuous, i.e.,

$$||g_{w_1} - g_{w_2}||_{\infty} \le J ||w_1 - w_2||. \tag{24}$$

**Theorem 5** (Infinite Rademacher Bound). Under assumptions (H1) and (H3), suppose that  $S^*(g) > S^* > 0$  and  $L(g) \leq L$  for all  $g \in \mathcal{G}$ . Furthermore, assume that  $p \geq n$ . Then, for any covering precision  $\tilde{\varepsilon} > 0$ ,

$$\mathcal{R}_{n,\mu}(\mathcal{F}) \qquad 389$$

$$\leq K \max\left\{\sqrt{\frac{1}{n}}, \frac{L}{S^*}\sqrt{\frac{cp}{nd}}\sqrt{\log(1+60WJ\tilde{\varepsilon}^{-1})}, \quad 390\right\}$$

$$2\exp\left(-\frac{dS^{*2}}{8cL^2}\right), \frac{J}{S^*}\tilde{\varepsilon}\right\}. \quad (25) \quad 391$$

where K > 0 is an absolute constant independent 392 of  $p, n, d, S^*, c, L, J, \tilde{\varepsilon}, W$ . 393

Proof. Given any discontinuous classifier  $f_w = \sup_{g_w \in \mathcal{G}} g_w \in \mathcal{G}$ , define its Lipschitz continuous approximation for  $\gamma > 0$  as

$$F_{f_w} = \operatorname{sgn}_{\gamma} \circ g_w, \tag{26}$$

where 398

$$\operatorname{sgn}_{\gamma}(t) := \begin{cases} -1, & t \leq -\gamma, \\ \frac{t}{\gamma}, & t \in [-\gamma, \gamma], \\ 1, & t \geq \gamma. \end{cases} \tag{27}$$

This approximation satisfies the useful property that both  $F_{f_w}$  and the absolute difference  $|f_w - F_{f_w}|$  are Lipschitz continuous in both the input space  $\mathcal{X}$  and the weight space  $\mathcal{W}$ , with

$$L(|\operatorname{sgn}_{\gamma} \circ g_w - \operatorname{sgn} \circ g_w|) = L(\operatorname{sgn}_{\gamma} \circ g_w) = \frac{L(g_w)}{\gamma}.$$
(28) 404

Using Lipschitz-continuous surrogates  $F_f$ , we decompose the Rademacher complexity into a smooth component, to which the analysis of Bubeck and Sellke [9] applies, and a residual term.

$$\mathcal{R}_{n,\mu}(\mathcal{F}) = \frac{1}{n} \mathbb{E}_{\sigma,x} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} \sigma_i f(x_i) \right| \right]$$
 409

$$\leq \frac{1}{n} \mathbb{E}_{\sigma, x} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} \sigma_i F_f(x_i) \right| \right]$$
 410

$$+ \frac{1}{n} \mathbb{E}_{\sigma, x} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} \sigma_i (f - F_f)(x_i) \right| \right]$$
 411

$$\leq C_1 \frac{1}{\sqrt{n}} + C_2 \frac{L}{\gamma} \sqrt{\frac{cp}{nd}} \sqrt{\log(1 + 60WJ/\tilde{\varepsilon})}$$
 412

$$+ \frac{1}{n} \mathbb{E}_{\sigma,x} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} \sigma_i (f - F_f)(x_i) \right| \right]. \quad (29) \quad 41$$

Here the parameter  $\tilde{\varepsilon} > 0$  is related to a  $\tilde{\varepsilon}$ -net 414 of  $\mathcal{W}$ , which we denote by  $\mathcal{W}_{\tilde{\varepsilon}}$ . Note, that  $|\mathcal{W}_{\tilde{\varepsilon}}| \leq 415$   $(1+60WJ\tilde{\varepsilon}^{-1})^p$  (see e.g. [11] Corollary 4.2.13) so the same holds true for the induced net  $\mathcal{F}_{\tilde{\varepsilon}} = \{\operatorname{sgn} \circ g_w : w \in \mathcal{W}_{\tilde{\varepsilon}}\}$ , which allows us to treat the remaining expectation by subdividing the supremum: 419

420 
$$\frac{1}{n} \mathbb{E}_{\sigma,x} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} \sigma_i (f - F_f)(x_i) \right| \right]$$

$$= \frac{1}{n} \mathbb{E}_{\sigma,x} \left[ \sup_{\substack{w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}} \\ \|w - w_{\tilde{\varepsilon}}\| \leq \tilde{\varepsilon}}} \left| \sum_{i=1}^{n} \sigma_i (f_w - F_{f_w})(x_i) \right| \right]$$

$$\leq \frac{1}{n} \mathbb{E}_x \Big[ \sup_{w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}}} \sum_{i=1}^n |f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}|(x_i) \Big]$$

$$+ \frac{1}{n} \mathbb{E}_{x} \left[ \sup_{\substack{w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}} \\ \|w - w_{\tilde{\varepsilon}}\| \leq \tilde{\varepsilon}}} \sum_{i=1}^{n} \left| |f_{w} - F_{f_{w}}| \right| \right]$$

$$-|f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}| |(x_i)|. \quad (30)$$

By Lipschitz continuity of the parameterization and of  $|f - F_f|$  (Equation 28), we obtain

$$||f_w - F_{f_w}| - |f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}||_{\infty} \le \frac{J}{\gamma} \tilde{\varepsilon}$$
 (31)

for any  $w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}}$  and  $w \in B_{\tilde{\varepsilon}}(w_{\tilde{\varepsilon}})$ , so that

429 
$$\frac{1}{n} \mathbb{E}_x \left[ \sup_{w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}}} \sum_{i=1}^n \left| |f_w - F_{f_w}|(x_i) \right| \right]$$

$$-|f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}|(x_i)| \bigg] \leq \frac{J}{\gamma} \tilde{\varepsilon}. \quad (32)$$

Note, that the expectation of the maximum of N subgaussians  $X_1, \ldots, X_N$  with variance proxy  $\sigma^2$  scales as

$$\mathbb{E}\left[\max_{1\leq i\leq N}|X_i|\right]\leq \sigma\sqrt{2\log\left(2N\right)},\tag{33}$$

see for instance [13]. The first expectation in Equation 30 can be bounded using Equation 33, since it corresponds – up to centering – to a maximum of sub-Gaussian random variables with variance proxy  $\sigma^2 = \frac{L^2}{\gamma^2} \frac{cn}{d}$ . Therefore,

$$\frac{1}{n} \mathbb{E}_x \Big[ \sup_{w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}}} \sum_{i=1}^n |f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}|(x_i) \Big] =$$

441 
$$\frac{1}{n} \mathbb{E}_x \left[ \sup_{w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}}} \sum_{i=1}^n |f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}|(x_i) - \mathbb{E}[|f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}|] \right]$$

$$+ \sup_{w_{\bar{\varepsilon}} \in \mathcal{W}_{\bar{\varepsilon}}} \mathbb{E}[|f_{w_{\bar{\varepsilon}}} - F_{f_{w_{\bar{\varepsilon}}}}|]$$

$$\leq C_3 \frac{L}{\gamma} \sqrt{\frac{c}{nd}} \sqrt{p \log(1 + 60WJ\tilde{\varepsilon}^{-1})}$$

$$+ \sup_{w_{\tilde{\varepsilon}} \in \mathcal{W}_{\tilde{\varepsilon}}} \mathbb{E}_x[|f_{w_{\tilde{\varepsilon}}} - F_{f_{w_{\tilde{\varepsilon}}}}|]. \quad (34)$$

Finally, for every  $f \in \mathcal{F}$ .

446 
$$\mathbb{E}_{x}[|f - F_{f}|] = \int_{\mathcal{X}} |f(x) - F_{f}(x)| \, d\mu(x)$$

$$\leq \mathbb{P}(g(x) \in [-\gamma, \gamma]). \quad (35)$$

Choosing  $\gamma = \frac{S^*(g)}{2}$ , we obtain by the definitions of co-margin, and once again isoperimetry (since

the co-margin inherits the Lipschitzness from g by design) 450

$$\mathbb{P}\left(g(x) \in [-\gamma, \gamma]\right) = \mathbb{P}\left(|g(x)| \le \frac{S^*(g)}{2}\right)$$
 452

$$\leq \mathbb{P}\left(|h_g^*(x) - S^*(g)| \geq \frac{S^*(g)}{2}\right) \tag{453}$$

$$\leq 2 \exp\left(-\frac{d \, S^*(g)^2}{8cL(g)^2}\right) \leq 2 \exp\left(-\frac{d \, S^{*2}}{8cL^2}\right)$$
 454

$$= 2 \exp \left(-\frac{d\,\bar{S}^{*2}}{8c}\right). \quad (36) \quad {}_{455}$$

Putting it all together, we have

$$\mathcal{R}_{n,\mu}(\mathcal{F})$$
 457

$$\leq C_1 \frac{1}{\sqrt{n}} + C_2' \frac{L}{S^*} \sqrt{\frac{c}{nd}} \sqrt{p \log(1 + 60WJ\tilde{\varepsilon}^{-1})} + \frac{2J}{S^*} \tilde{\varepsilon} \quad \text{456}$$

$$+ \, 2 \exp \biggl( - \frac{d \, S^{* \, 2}}{8 c L^2} \biggr) \,, \quad (37) \quad {\it 45}$$

for absolute constants  $C_1, C_2'$ , independent of  $p, n, d, S^*, c, L, J, \tilde{\varepsilon}, W$ .

## A.3 Proof of the Law of Robustness 462

Next, we provide the proof of the law of robustness description for classification problems.

**Corollary 6** (Law of Robustness for Discontinuous Functions). Assume we are in the setting of Theorem 4.66 Let  $p := \log |\mathcal{F}| \geq n$ . Fix  $\varepsilon, \delta \in (0,1)$  and 467 consider the 0-1 loss  $\ell_{0-1}$ . There exists an absolute 468 constant K > 0 such that, if 469

1. the minimal risk 
$$\sigma^2 := \min_{f \in \mathcal{F}} R_{0-1}(f)$$
 satisfies  $\sigma^2 \geq \varepsilon$ , and 471

2. the sample size n is large enough to ensure (i) 472

$$\frac{K}{\sqrt{n}} < \frac{\varepsilon}{3} \text{ and (ii) } \sqrt{\frac{2\log(2/\delta)}{n}} < \frac{\varepsilon}{2},$$

then with probability at least  $1 - \delta$  (over the sample), 474 the following holds uniformly for all  $f \in \mathcal{F}$ : 475

$$\hat{R}_{0-1}(f) \le \sigma^2 - \varepsilon \quad \Longrightarrow \tag{476}$$

$$S(f) < \max \left\{ \frac{3K}{\varepsilon} \sqrt{\frac{c \log |\mathcal{F}|}{nd}}, \sqrt{\frac{8c}{d} \log \left(\frac{6K}{\varepsilon}\right)} \right\}. \tag{38}$$

*Proof.* Let K > 0 be an absolute constant such that 476 Equation 22 holds, and define the threshold stability 478

$$S_* = S_*(p, n, d, \varepsilon) \tag{480}$$

$$:= \max \left\{ \frac{3K}{\varepsilon} \sqrt{\frac{c \log |\mathcal{F}|}{nd}}, \sqrt{\frac{8c}{d} \log \left(\frac{6K}{\varepsilon}\right)} \right\}. \tag{39}$$

525

526

528

539

542

Then, Theorem 4, together with condition 2(i), implies that 483

484 
$$\mathcal{R}_{n,\mu}(\mathcal{F}_{S_*})$$
485  $\leq K \max \left\{ \frac{1}{\sqrt{n}}, \frac{\sqrt{c}}{S_*} \sqrt{\frac{\log |\mathcal{F}|}{nd}}, 2 \exp\left(-\frac{dS_*^2}{8c}\right) \right\}$ 
486  $\leq \varepsilon/3, \quad (40)$ 

where  $\mathcal{F}_{S_*} := \{ f \in \mathcal{F} : S(f) \geq S_* \}$  is the subset 487 of functions in  $\mathcal{F}$  with stability at least  $S_*$ . Hence, 488 applying the standard generalization inequality in 489 terms of the Rademacher complexity [2], together 490 with condition 2(ii), yields that with probability at least  $1 - \delta$ :

$$\sup_{f \in \mathcal{F}_{S_*}} \left( R_{0-1}(f) - \hat{R}_{0-1}(f) \right)$$

494 
$$\leq 2\mathcal{R}_{n,\mu}(\ell_{0-1} \circ \mathcal{F}_{S_*}) + \sqrt{\frac{2\log(2/\delta)}{n}}$$
495 
$$\leq \mathcal{R}_{n,\mu}(\mathcal{F}_{S_*}) + \frac{\varepsilon}{2} < \varepsilon, \quad (41)$$

where we additionally used 496

$$\mathcal{R}_{n,\mu}(\ell_{0-1} \circ \mathcal{F}) \le \frac{1}{2} \mathcal{R}_{n,\mu}(\mathcal{F}), \tag{42}$$

in the second step. In particular, we can bound the 498 probability 499

500 
$$\mathbb{P}(\forall f \in \mathcal{F}_{S_*} : \hat{R}_{0-1}(f) > \sigma^2 - \varepsilon)$$
  
501  $\geq \mathbb{P}(\forall f \in \mathcal{F}_{S_*} : R_{0-1}(f) - \hat{R}_{0-1}(f) < \varepsilon) \geq 1 - \delta,$ 
(43)

where the first inequality follows from 502

503 
$$R_{0-1}(f) - \hat{R}_{0-1}(f) < \varepsilon \stackrel{\text{condition 1.}}{\Longrightarrow}$$
504  $\sigma^2 - \hat{R}_{0-1}(f) < \varepsilon \implies \hat{R}_{0-1}(f) > \sigma^2 - \varepsilon.$  (44)

505 events 506

507 
$$1 - \delta \leq \mathbb{P}(\forall f \in \mathcal{F}_{S_*} : \hat{R}_{0-1}(f) > \sigma^2 - \varepsilon)$$
508 
$$= \mathbb{P}(\forall f \in \mathcal{F} : \hat{R}_{0-1}(f) > \sigma^2 - \varepsilon)$$
509 
$$+ \mathbb{P}(\exists f \in \mathcal{F}_{S_*}^c : \hat{R}_{0-1}(f) \leq \sigma^2 - \varepsilon), \quad (45)$$

enables us to easily recognize that the expression exactly characterizes the probability that the following 511 implication, and thereby the result, holds uniformly 512 for all  $f \in \mathcal{F}$ : 513

$$\hat{R}_{0-1}(f) \le \sigma^2 - \varepsilon \implies S(f) < S_*. \tag{46}$$

Indeed, the implication above holds if, for a given 515 data sample  $(x_i, y_i)_{i=1}^n$ , either 516

514

• no function 
$$f \in \mathcal{F}$$
 satisfies  $\hat{R}_{0-1}(f) \leq \sigma^2 - \varepsilon$ ,
or

• any such f lies in  $\mathcal{F}_{S_a}^c$ , that is,  $S(f) < S_*$ , which is the case with probability at least  $1 - \delta$  due 520 to Equation 45. 521

With the same reasoning and Theorem 5, we obtain a law of robustness for infinite classes. 523

Corollary 7 (Law of Robustness for Infinite Function Classes). Assume we are in the setting of Theorem 5, and fix  $\varepsilon, \delta \in (0,1)$ . Consider the 0-1 loss  $\ell_{0-1}$ . There exists an absolute constant K>0 such that, if

- 1. the minimal risk  $\sigma^2 := \min_{f \in \mathcal{F}} R_{0-1}(f)$  satisfies 529 530
- 2. the sample size n is large enough so that (i) 531  $\frac{K}{\sqrt{n}} < \frac{\varepsilon}{3}$  and (ii)  $\sqrt{\frac{2\log(2/\delta)}{n}} < \frac{\varepsilon}{2}$ 532

then with probability at least  $1 - \delta$ , for all  $\tilde{\varepsilon} > 0$ , 533 the following holds uniformly for all  $g \in \mathcal{G}$  and  $f_q = sgn \circ g$ : 535

$$\hat{R}_{\theta-1}(f_g) \le \sigma^2 - \varepsilon \implies 536$$

$$\frac{S^*(g)}{L(g)} < \max \left\{ \frac{3K}{\varepsilon} \sqrt{\frac{p}{nd}} \sqrt{c \log(1 + 60WJ\tilde{\varepsilon}^{-1})} \right\} 537$$

$$, \sqrt{\frac{8c}{d} \log\left(\frac{6K}{\varepsilon}\right)} \right\}. (47) 538$$

## **Multi-Class Classification**

In this section, we briefly outline how our results ex-540 tend to categorical distributions with  $C \in \mathbb{N}$  classes. We assume that a classifier is given by

$$f: \mathcal{X} \to \{0, 1\}^{\mathcal{C}},$$
 (48) 543

with exactly one non-zero entry for each  $x \in \mathcal{X}$ . The adaptations of the conditions in (H3) to the 545  $\sigma^2 - \hat{R}_{0-1}(f) < \varepsilon \implies \hat{R}_{0-1}(f) > \sigma^2 - \varepsilon.$  (44) multi-class setting can be formalized as follows.

Decomposing this probability into two disjoint (H3)' The hypothesis class has the form  $\mathcal{F}$ 547  $\operatorname{argmax} \circ \mathcal{G}$ , where  $\mathcal{G} = \{g_w : \mathcal{X} \to [0,1]^{\mathcal{C}} : w \in \mathcal{G}\}$ W is a parameterized family of Lipschitz func-549 tions. The parameter space  $\mathcal{W} \subset \mathbb{R}^p$  is bounded with  $diam(W) \leq W$ , and the parameterization 551 is Lipschitz:

$$\|g_{w_1} - g_{w_2}\|_{\infty} \le J \|w_1 - w_2\|.$$
 (49) 553

Thus, we can interpret  $g \in \mathcal{G}$  as representing the class probabilities. 555

**Remark 8.** For binary classification, i.e. C = 2, the 556 classifiers are of the form  $f: \mathcal{X} \to \{0,1\}^2$ , instead of  $f: \mathcal{X} \to \{-1, 1\}$ , as considered earlier. However, 558 one can translate between these representations by 559 post-composing with either

$$\alpha(x_1, x_2) := x_1 - x_2 \quad or \quad \beta(x) := \left(\frac{x+1}{2}, \frac{1-x}{2}\right).$$
(50) 5

603

612

626

628

630

631

By the contraction principle for Rademacher complexity, it is therefore sufficient to compute the com-563 plexity for one of these models. 564

As in the binary case, our proofs start by consid-565 ering the Rademacher complexity of the function 566 567

$$\mathcal{R}_{n,\mu}(\mathcal{F}) = \frac{1}{n} \mathbb{E}^{\sigma_{ij},x_i} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} \sum_{j=1}^{\mathcal{C}} \sigma_{ij} f_j(x_i) \right| \right]$$
(51)

$$\leq \sum_{j=1}^{\mathcal{C}} \frac{1}{n} \mathbb{E}^{\sigma_{ij}, x_i} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} \sigma_{ij} f_j(x_i) \right| \right].$$
(52)

Each summand corresponds to a binary classification 570 problem with a one-vs-all classifier  $f_j$ . Indeed,  $f_j$  is  $\frac{2}{S(f)-t}$ -Lipschitz on  $A_t(f)$ . Transforming via

$$f_j \mapsto 2f_j - 1 : \mathcal{X} \to \{-1, 1\},$$
 (53)

we can follow the same reasoning as in A.2, obtaining, 574 up to a linear factor of C, the same result as the first 575 part of Theorem 4, generalized to the multi-class 576 setting.

Similarly, under assumption (H3), we can write 578

$$2f_j - 1 = \operatorname{sgn}(g_j - \max_{i \neq j} g_i(x)), \tag{54}$$

which allows us to proceed as in Theorem 5 to obtain 580 a multi-class generalization of Theorem 5 and Corollary 7. The only minor difference lies in bounding the term in Equation 35: 583

$$\mathbb{E}[|f_j - F_{f_j}|] \le \mathbb{P}[|g_j(x) - \max_{i \ne j} g_i(x)| \le \gamma]. \quad (55)$$

Choosing  $\gamma = \frac{S^*(g)}{2}$ , we use that for all j,  $|g_j(x) - \max_{i \neq j} g_i(x)| > h_g^*(x)$ , which yields 585 586

$$\mathbb{P}[|g_j(x) - \max_{i \neq j} g_i(x)| \le \frac{S^*(g)}{2}]$$
 (56)

$$\leq \mathbb{P}[|h_g^*(x) - S^*(f)| \geq \frac{S^*(g)}{2}] \qquad (57)$$

$$\leq 2\exp\left(-\frac{dS^*(g)^2}{8cL(g)^2}\right) \tag{58}$$

$$\leq 2\exp\left(-\frac{dS^{*2}}{8cL^2}\right) \tag{59}$$

$$= 2\exp\left(-\frac{d\bar{S}^{*2}}{8c}\right). \tag{60}$$

We conclude that all results extend naturally to 592 the multi-class case. The main concepts are summa-593 rized below.

### • Isoperimetry:

595

596

$$\mathbb{P}(\|f(x) - \mathbb{E}[f]\|_{\infty} \ge t) \le 2\exp\left(-\frac{dt^2}{2cL^2}\right) \tag{61}$$

Rademacher complexity

$$\mathcal{R}_{n,\mu}(\mathcal{F}) = \frac{1}{n} \mathbb{E}^{\sigma_{i,j}, x_i} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sum_{j=1}^{\mathcal{C}} \sigma_{ij} f_j(x_i) \right| \right]$$
(62) 598

• Margin

$$h_f(x) = \sum_{j=1}^{C} h_f^j(x)$$
 (63) 600

$$h_f^j(x) := \inf\{\|x - z\|_2 : f(z) \neq j, z \in \mathbb{R}^d\}$$
 (64) 602

• Class stability

$$S(f) = \sum_{j=1}^{\mathcal{C}} S(f)^j, \quad S(f)^j := \mathbb{E}[h_f^j]$$
 (65) 60

• Co-margin 605

$$h_g^*(x) = \sum_{j=1}^{C} h_g^{*j}(x)$$
 (66) 606

$$h_g^{*j}(x) := \max\left(0, g_j(x) - \max_{i \neq j} g_i(x)\right)$$
 (67) 608

 Co-stability 609

$$S^*(g) = \sum_{j=1}^{\mathcal{C}} S^{*j}(g), \quad S^{*j}(g) := \mathbb{E}[h_g^{*j}] \quad (68)$$

### Experimental Details for Stabil- 611 A.5ity Measurement

Training setup. To empirically validate our ro- 613 bustness law, we trained fully connected MLPs on MNIST and CIFAR-10 datasets. Each model has 4 615 hidden layers with widths  $w \in \{128, 256, 512,$ 616 1024, 2048 for MNIST and up to w = 1024 for CI-617 FAR10. All models use ReLU activations, batch normalization, and were initialized with standard parametrization. Training was conducted using the Adam optimizer [14] for the embedding and output layers, and the Muon optimizer [15] for the hidden layers. Models were trained with a batch size of 256 and learning rate  $10^{-3}$ , until at least 99% training accuracy was achieved, ensuring (near) interpolation. We further used sharpness-aware optimization based on [16, 17] to reduce variance of the normalized 627 co-stability on MNIST.

Parameter counts and normalization. For each model, we recorded the total number of trainable parameters p, input dimension d, and total number of training samples n.

Stability estimation. Class stability S(f) was computed using adversarial perturbation analysis. We performed a suite of  $\ell_2$ -based attacks (FGSM, PGD, DeepFool, and L2PGD [18–20]) using the Foolbox library [21]. For each input x, we recorded the minimum perturbation norm required to change the classifier's prediction, over a grid of radii  $\mathbf{r} = (0.002, 0.01, 0.05, 0.1)$ . The final stability score S(f) was taken as the average  $\ell_2$  distance across the dataset.

Normalized Co-Stability estimation. The empirical co-stability  $S^*(g)$  is computed via the multiclass margin

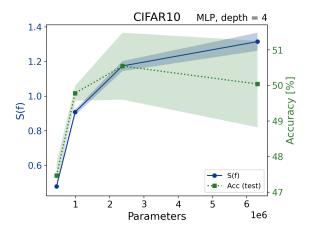
$$g_j(x) - \max_{i \neq j} g_i(x), \qquad j = \arg\max_i g_i(x), \quad (69)$$

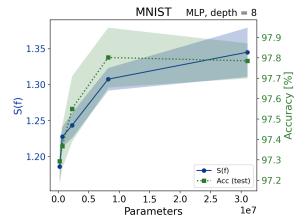
averaged over the dataset. We estimate the Lipschitz constant L(g) using the efficient ECLIPSE method [12], and report the normalized ratio  $S^*(g)/L(g)$  as a function of model size.

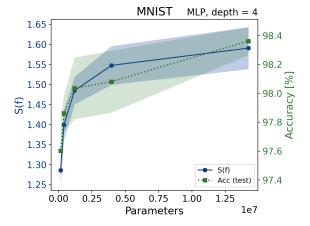
Implementation. Training and evaluation code is implemented in PyTorch [22]. For MLPs, images were flattened to vectors. Attack evaluations were conducted over the full dataset (train and test).

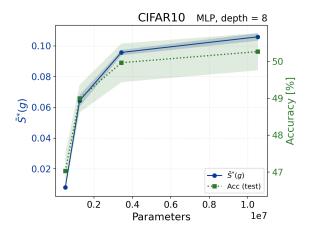
**Reproducibility.** All experiments were run with multiple random seeds  $\{0, 1, 2, 3, 4\}$ , and mean with standard deviation are reported. Our code is available at anonymous GitHub.

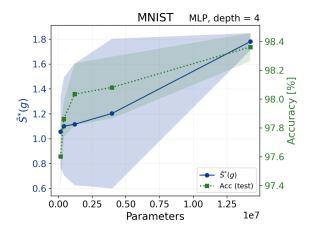
Results. Figure A.1 shows that, for MLPs, both class stability S(f) and normalized co-stability  $S^*(g)/L(g)$  increase consistently with model size. The observed saturation of (normalized co-) stability aligns with theoretical intuition: the Bayes classifier admits a finite (normalized co-) stability level, and pushing beyond this level necessarily reduces accuracy – an instance of the robustness/accuracy trade-off extensively discussed in the literature [23–25]. Accordingly, we expect stability to plateau once models approach the Bayes decision boundary. For CIFAR-10, although test accuracy remains far below the Bayes optimal (around 50%), the same reasoning applies relative to the best classifier achievable within the restricted MLP architecture.

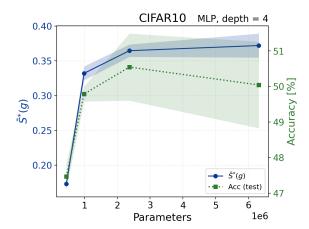


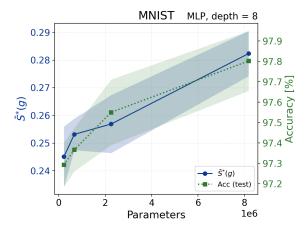












**Figure A.1.** Stability measures for 4- and 8-layer MLPs trained on MNIST and CIFAR-10.