046

047

052

053

054

000

On the Geometry of Regularization in Adversarial Training: High-Dimensional Asymptotics and Generalization Bounds

Anonymous Authors¹

Abstract

Regularization, whether explicit in terms of a penalty in the loss or implicit in the choice of algorithm, is a cornerstone of modern machine learning. Indeed, controlling the complexity of the model class is particularly important when data is scarce, noisy or contaminated, as it translates a statistical belief on the underlying structure of the data. This work investigates the question of how to choose the regularization norm $\|\cdot\|$ in the context of high-dimensional adversarial training for binary classification. To this end, we first derive an exact asymptotic description of the robust, regularized empirical risk minimizer for various types of adversarial attacks and regularization norms (including non- ℓ_p norms). We complement this analysis with a uniform convergence analysis, deriving bounds on the Rademacher Complexity for this class of problems. Leveraging our theoretical results, we quantitatively characterize the relationship between perturbation size and the optimal choice of $\|\cdot\|$, confirming the intuition that, in the data scarce regime, the type of regularization becomes increasingly important for adversarial training as perturbations grow in size.

1. Introduction

Despite all its successes, deep learning still underperforms spectacularly in worst-case situations, when models face innocent-looking data which are adversarially crafted for eliciting erroneous or undesired outputs. Since the discovery of these failure modes in computer vision (Szegedy et al., 2014) and their re-discovery, more recently, in other modalities including text (Zou et al., 2023), considerable effort has been put in designing algorithms for training models which are robust against these adversarial attacks. In the context of supervised learning problems, a principled approach consists of appropriately modifying standard empirical risk minimization: a parametric model is fit by minimizing a worst-case empirical risk, where "worst-case" refers to an assumed threat model. For example, in computer vision, a threat model of ℓ_{∞} perturbations translates the assumption that images whose pixels only differ by a little should share the same label. Despite its conceptual clarity and proven ability to return robust models, a major drawback of this method, known as robust empirical risk minimization (RERM) or adversarial training (Goodfellow et al., 2015; Madry et al., 2018), is that it often comes with a performance tradeoff, besides being computationally more intensive than standard ERM. Indeed, it has been observed that model accuracy is often compromised for better robustness (Tsipras et al., 2019; Zhang et al., 2019). To make matters worse, neural networks often exhibit a large gap between their robust train and test performances in standard computer vision benchmarks (Rice et al., 2020).

Many empirical efforts in addressing these statistical limitations of RERM have focused on either increasing the amount of labeled (Wang et al., 2023) or unlabeled (Carmon et al., 2019; Zhai et al., 2019) data, or on painstakingly re-imagining several of the design choices of deep learning (such at the loss function (Zhai et al., 2019), model averaging (Chen et al., 2021; Rebuffi et al., 2021) and more). Despite the apparent empirical challenges, simple guidelines on how different choices affect the statistical efficiency of RERM are clearly missing, even in simple models.

In this work, we make a step towards theoretically filling this gap by investigating model selection in RERM, and how it relates to robust and standard generalization error. In particular, we focus on the oldest model selection method: (weight) *regularization*. Following a large body of work originating in high-dimensional statistics (Krogh & Hertz, 1991; Seung et al., 1992; Bean et al., 2013a; Thrampoulidis et al., 2018; Aubin et al., 2020; Vilucchio et al., 2024), we study this fundamental question *asymptotically*, when both the input dimension and the number of training samples grow to infinity while keeping their ratio constant, and under a Gaussian setting. While it is customary in this literature to study which *values* of regularization coefficients

 ¹Anonymous Institution, Anonymous City, Anonymous Region,
 Anonymous Country. Correspondence to: Anonymous Author
 <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

vield the best test errors (balancing empirical fitness with model complexity), we instead analyze the optimality of a 057 type of regularization. A motivation for this comes from a 058 separate line of work in uniform convergence bounds that 059 stresses the importance of the type of regularization for ro-060 bust generalization (Yin et al., 2019; Awasthi et al., 2020; 061 Tsilivis et al., 2024). Borrowing from this line of work, 062 which mainly offers qualitative bounds, and reinforcing it 063 with new findings, we demonstrate, via sharp asymptotic de-064 scriptions of the errors in (regularized) RERM for a variety 065 of different perturbation and regularization norms, that reg-066 ularization becomes increasingly important in RERM as the 067 perturbation strength grows in size. This allows us to get an 068 exact description of the relationship between optimal type 069 of regularization and strength of perturbation, and discuss 070 how regularization affects the tradeoff between robustness 071 and accuracy.

To summarize, our main contributions in this work are the following:

074

081

088

089

090

091

092

093

094

095

096

097

- 075 1. We derive an exact asymptotic description of the per-076 formance of reguralized RERMs for a variety of per-077 turbation and regularization norms. In addition to the 078 usually studied ℓ_p , we consider $\|\cdot\|_{\Sigma}$ norms (induced 079 by a positive symmetric matrix Σ), which allow us to separate the effect of a perturbation on different features of the input. 082
- 083 2. We show uniform convergence bounds for this class of problems (i.e., $\|\cdot\|_{\Sigma}$ regularized), by establishing new results on the Rademacher complexities for several classes of linear hypothesis classes under adversarial perturbations. 087
 - 3. Leveraging the theoretical results above, we show that regularizing with the *dual* norm of the perturbation can yield significant benefits in terms of robustness and accuracy, compared to other regularization choices. In particular, our analysis permits a precise characterization of the relationship between the perturbation geometry and the optimal type of regularization. It further allows a decomposition of the contribution of regularization in terms of standard and robust (test) error.

098 Our results can be seen as positive news. Indeed, the main 099 implication of our work for robust machine learning practice 100 is that model selection, in the form of either explicit or implicit regularization, plays a more important role in robust ERM than in standard ERM. In the context of robust deep learning practice, model selection is often implicit in the 104 choice of architecture, learning algorithm, stopping time, 105 hyperparameters, etc. Our theoretical analysis in the context 106 of simple adversarial tasks highlights the importance of these choices, as they can be crucial to the outcome in terms of robustness and performance. 109

Finally, while typical-case and worst-case analyses usually appear as opposites in the statistical learning literature, we believe our work nicely illustrates how these two approaches to studying generalization can be combined in a complementary way to yield precise answers with both explanatory and predictive powers.

1.1. Further related Work

Exact asymptotics: Our results on the exact asymptotics of adversarial training build on an established body of literature that spans high-dimensional probability (Thrampoulidis et al., 2014; 2015; Taheri et al., 2023), statistical physics (Mignacco et al., 2020; Gerace et al., 2021; Bordelon et al., 2020; Loureiro et al., 2021; Okajima et al., 2023; Adomaityte et al., 2024; 2023) and random matrix theory (Bean et al., 2013b; Mai et al., 2019; Liao et al., 2020; Mei & Montanari, 2022; Xiao et al., 2022; Schröder et al., 2023). Our study is also motivated by recent efforts to understand Gaussian universality (Goldt et al., 2021; Montanari & Saeed, 2022; Dandi et al., 2023). These works suggest that simple models for the covariates can have a broader scope in the context of high-dimensional generalized linear estimation, often mirroring real-world datasets. From a technical perspective, this phenomenon arises due to strong concentration properties in the high-dimensional regime, which imply some universality properties of the generalization error with respect to the covariate distribution (Tao & Vu, 2010; Donoho & Tanner, 2009; Wei et al., 2022; Dudeja et al., 2023).

Adversarial training: Robust empirical risk minimization, i.e. adversarial training, was first introduced for computer vision applications (Goodfellow et al., 2015; Madry et al., 2018). A large body of work is devoted to the study of applied methods for improving its computational (Shafahi et al., 2019; Rice et al., 2020) and statistical (Zhai et al., 2019; Chen et al., 2021; Wang et al., 2023) properties. Theoretically, robust training has been considered before in both the case of Gaussian mixture models (Bhagoji et al., 2019; Dan et al., 2020; Javanmard & Soltanolkotabi, 2022) and linear regression (Raghunathan et al., 2020; Taheri et al., 2023; Dohmatob & Scetbon, 2024). Of particular interest is the work of Tanner et al. (2024), which recently derived high dimensional asymptotics for binary classification with ℓ_2 regularization, considering perturbations in a general $\|\cdot\|_{\Sigma}$ norm. In our work, we study the effect of regularization, providing exact asymptotics for any combination of ℓ_p perturbation and regularization norm, while extending (Tanner et al., 2024) for various $\|\cdot\|_A$ regularization norms (where A is a positive symmetric matrix).

Statistical learning theory: The role of regularization in statistical inference traces back to the work of Tikhonov

(1963) and plays a central role in statistical learning the-111 ory, directly inspiring general inductive principles such as 112 Structural Risk Minimization (Vapnik, 1998) and practical 113 methods that realize this principle, such as SVMs (Cortes & 114 Vapnik, 1995). Uniform convergence bounds, quantities that 115 upper bound the difference between empirical and expected 116 risk of any predictor *uniformly* inside a hypothesis class, 117 were originally stated as a function of the VC dimension of 118 the class (Vapnik & Chervonenkis, 1971). The Rademacher 119 complexity of the class (Koltchinskii, 2001) is known to 120 typically provide finer guarantees (Kakade et al., 2008). 121 Several recent papers derive such results in the context of 122 adversarially robust classification for linear predictors and neural networks (Yin et al., 2019; Awasthi et al., 2020; Xiao 124 et al., 2024). Based on these results, (Tsilivis et al., 2024) 125 highlighted the importance of the (implicit) regularization 126 in RERM with linear models, by showing the effect of the 127 learning algorithm and the architecture on the robustness 128 of the final predictor. In our work, we consider, instead, 129 explicit regularization and more general perturbation (and 130 regularization) geometries. 131

2. Setting Specification

132

133

138

157

158 159

160

161

162

163

164

134 We consider a binary classification task with training data 135 $S = \{(x_i, y_i)\}_{i=1}^n$, where $x_i \in \mathbb{R}^d$ and $y_i \in \{-1, +1\}$ 136 are sampled independently from a distribution \mathcal{D} of the 137 following form:

$$P(\boldsymbol{x}, y) = \int_{\mathbb{R}^d} \mathrm{d}\boldsymbol{w}_{\star} \mathbb{P}_{\mathrm{out}}\left(y \Big| \frac{\langle \boldsymbol{w}_{\star}, \boldsymbol{x} \rangle}{\sqrt{d}}\right) P_{\mathrm{in}}(\boldsymbol{x}) P_{\boldsymbol{w}}(\boldsymbol{w}_{\star}),$$

$$(1)$$

142 where P_{in} is a probability density function over \mathbb{R}^d and 143 $\mathbb{P}_{out}: \mathbb{R} \to [0,1]$ encodes our assumption that the label 144 is a (potentially non-deterministic) linear function of the 145 input x with teacher weights $w_{\star} \in \mathbb{R}^{d}$. For example, a 146 noiseless problem corresponds to $\mathbb{P}_{out}(y|z) = \delta(y-z)$, 147 while we can incorporate noise by using the *probit* model: 148 $\mathbb{P}_{\mathrm{out}}(y|z) = 1/2 \operatorname{erfc}(-yz/\sqrt{2}\tau)$, where $\tau > 0$ controls the 149 label noise. We assume that $w_{\star} \in \mathbb{R}^d$ is drawn from a prior 150 distribution $P_{\rm w}$. 151

152 Given the training data S, our objective is to investigate 153 the robustness and accuracy of linear classifiers $\hat{y}(\hat{w}, x) =$ 154 sign $(\langle \hat{w}, x \rangle / \sqrt{d})$, where $\hat{w} = \hat{w}(S)$ are learned from the 155 training data.

We define the robust generalization error as

$$E_{\rm rob}(\hat{\boldsymbol{w}}) = \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \left[\max_{\|\boldsymbol{\delta}\| \leq \varepsilon} \mathbb{1}(y \neq \hat{y}(\hat{\boldsymbol{w}}, \boldsymbol{x} + \boldsymbol{\delta})) \right], \quad (2)$$

where the pair (x, y) comes from the same distribution as the training data, and ε bounds the magnitude of adversarial perturbations under a specific choice of norm. The (standard) *generalization error* is defined as the rate of misclassification of the learnt predictor

$$E_{\text{gen}}(\hat{\boldsymbol{w}}) = \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}}[\mathbb{1}(y \neq \hat{y}(\hat{\boldsymbol{w}}, \boldsymbol{x}))].$$
(3)

Notice that for $\varepsilon = 0$: $E_{\rm rob}(\hat{\boldsymbol{w}}) = E_{\rm gen}(\hat{\boldsymbol{w}})$ for all $\hat{\boldsymbol{w}} \in \mathbb{R}^d$. We will frequently use the following decomposition of the robust generalization error into (standard) generalization error and *boundary error* $E_{\rm bnd}$:

$$E_{\rm rob}(\hat{\boldsymbol{w}}) = E_{\rm gen}(\hat{\boldsymbol{w}}) + E_{\rm bnd}(\hat{\boldsymbol{w}}), \qquad (4)$$

where E_{bnd} is defined as follows

$$E_{\text{bnd}}(\hat{\boldsymbol{w}}) = \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}} \bigg[\mathbb{1}(y = \hat{y}(\hat{\boldsymbol{w}}; \boldsymbol{x})) \max_{\|\boldsymbol{\delta}\| \leq \varepsilon} \mathbb{1}(y \neq \hat{y}(\hat{\boldsymbol{w}}, \boldsymbol{x} + \boldsymbol{\delta})) \bigg].$$
(5)

As its name suggests, $E_{\rm bnd}$ is the probability of a sample lying on (or near) the decision boundary, i.e., the probability that a sample is correctly classified without perturbation but incorrectly classified with it.

2.1. Robust Regularized Empirical Risk Minimization

Direct minimization of the robust generalization error of eq. (2) presents two main challenges: first, the objective function is non-convex due to the indicator function, and second, we only have access to a finite dataset rather than the full data-generating distribution. To address these issues, a widely adopted approach is to optimise the robust *empirical* (regularized) risk, defined as

$$\mathcal{L}(\boldsymbol{w}) = \sum_{i=1}^{n} \max_{\|\boldsymbol{\delta}_i\| \leq \varepsilon} g\left(y_i \frac{\langle \boldsymbol{w}, \boldsymbol{x}_i + \boldsymbol{\delta}_i \rangle}{\sqrt{d}}\right) + \lambda \widetilde{r}(\boldsymbol{w}), \quad (6)$$

where $g: \mathbb{R} \to \mathbb{R}_+$ is a non-increasing convex loss function that serves as a surrogate for the 0-1 loss, r(w) is a convex regularization term, and $\lambda \ge 0$ is a regularization parameter. The inner maximization over δ_i models the worst-case perturbation for each data point, constrained by the attack budget ε during training. Given the dataset S, we estimate the parameters of our model as

$$\hat{\boldsymbol{w}} \in \operatorname*{arg\,min}_{\boldsymbol{w} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{w}).$$
 (7)

The choice of loss function g, regularization r, and parameters ε and λ can significantly impact the model's accuracy and robustness.

In practice, eq. (7) is often solved with a first-order optimization method, such as gradient descent. Prior work (Soudry et al., 2018) has showed that optimizing the *unregularized* loss without any adversarial perturbations (eq. (6) for $\lambda, \varepsilon = 0$) with gradient descent is equivalent to eq. (6) with 165 the euclidean norm squared as a regularizer, where the free regularizer strength λ corresponds to the time duration of 167 the algorithm ($\lambda \rightarrow 0$ as the number of iterations goes to 168 ∞). Similar results can be obtained for different first-order algorithms (Gunasekar et al., 2018) (in particular, when 169 170 $r(w) = ||w||_p^p$, this corresponds to the family of steepest descent algorithms) as well as in the adversarial case 171 172 $(\varepsilon > 0)$ (Tsilivis et al., 2024). Therefore, studying eqs. (6) and (7) is equivalent to studying the solutions returned by a 174 first-order optimization algorithm.

3. Exact Asymptotics of Robust ERM

178 Our first technical result is an asymptotic description of 179 the properties of the solution of eqs. (6) and (7) in the pro-180 portional high-dimensional limit, under the assumption of 181 isotropic Gaussian distribution. While restrictive, this as-182 sumption is supported by recent theoretical advances show-183 ing that many learning problems exhibit universality: their 184 asymptotic behavior matches Gaussian predictions even 185 with non-Gaussian data (Goldt et al., 2022; Loureiro et al., 186 2021; Hu & Lu, 2023; Montanari & Saeed, 2022; Dandi 187 et al., 2023; Wei et al., 2022; Pesce et al., 2023; Gerace et al., 2024). While proving such correspondence for our 189 setting is outside the scope of this work, this suggests our 190 analysis of the Gaussian case can provide valuable insights 191 into practical adversarial training. 192

193 **3.1. Results for** ℓ_p **norms**

175

177

First, we consider the setting where the perturbations in eqs. (2) and (6) are constrained in their ℓ_p norm for $p \in$ (1, ∞]. More precisely, we make the following assumptions:

198 Assumption 3.1 (High-Dimensional Limit). We consider 199 the proportional high-dimensional regime where both the 200 number of training data n and input dimension d diverge to 201 infinity simultaneously at the same rate, while maintaining 202 a fixed ratio $\alpha := n/d$.

Assumption 3.2 (ℓ_p Norms). Let $\|\boldsymbol{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ denote the ℓ_p norm for $p \in (1, \infty]$, with p^* being its dual exponent $(1/p + 1/p^* = 1)$. The adversarial perturbations are constrained by an ℓ_p norm with parameter p, while for regularization we consider the function $\tilde{r}(\boldsymbol{w}) = \|\boldsymbol{w}\|_r^r$ where $r \in [1, \infty)$ is a parameter that can differ from p.

210 Assumption 3.3 (Scaling of Adversarial Norm Constraint). 211 We suppose that the value of ε scales with the dimension d212 such that $\varepsilon/d^{1/p^*} = O_d(1)$.

213 214 214 215 216 217 218 219 **Assumption 3.4** (Data Distribution). For each $i \in [n]$, the covariates $x_i \in \mathbb{R}^d$ are drawn i.i.d. from the data distribution $P_{in}(x) = \mathcal{N}_x(0, \mathrm{Id}_d)$. Then the corresponding y_i is sampled independently from the conditional distribution \mathbb{P}_{out} defined in eq. (1). The target weight vector $w_{\star} \in \mathbb{R}^d$ is drawn from a prior probability distribution P_w which is separable, i.e. $P_{\boldsymbol{w}}(\boldsymbol{w}) = \prod_{i=1}^{d} P_{\boldsymbol{w}}(w_i)$ for a distribution $P_{\boldsymbol{w}}$ in \mathbb{R} with finite variance $\operatorname{Var}(P_{\boldsymbol{w}}) = \rho < \infty$.

Under these assumptions, our first result states that in the high-dimensional limit, the robust generalization error associated with the RERM solution in eq. (6) asymptotically depends only on a few deterministic variables, known as the *summary statistics*, which can be computed by solving a set of low-dimensional self-consistent equations.

Theorem 3.5 (Limiting errors for ℓ_p norm). Let $\hat{w}(S) \in \mathbb{R}^d$ denote a solution of the RERM problem in eq. (6). Then, under Assumptions 3.1 to 3.4, the standard, robust and boundary generalization error of \hat{w} converge to the following deterministic quantities

$$E_{\text{gen}}(\hat{\boldsymbol{w}}) = \frac{1}{\pi} \operatorname{arccos}\left(m^{\star}/\sqrt{(\rho + \tau^{2})q^{\star}}\right),$$
$$E_{\text{bnd}}(\hat{\boldsymbol{w}}) = \int_{0}^{\varepsilon \frac{p^{\star}\sqrt{p^{\star}}}{\sqrt{q^{\star}}}} \operatorname{erfc}\left(\frac{-\frac{m^{\star}}{\sqrt{q^{\star}}}\nu}{\sqrt{2(\rho + \tau^{2} - m^{\star^{2}}/q^{\star})}}\right) \frac{e^{-\frac{\nu^{2}}{2}}}{\sqrt{2\pi}} \,\mathrm{d}\nu\,,$$
$$E_{\text{rob}}(\hat{\boldsymbol{w}}) = E_{\text{gen}}(\hat{\boldsymbol{w}}) + E_{\text{bnd}}(\hat{\boldsymbol{w}})$$

where m^*, q^*, P^* are the limiting values of the following summary statistics:

$$\frac{1}{d} \langle \boldsymbol{w}_{\star}, \hat{\boldsymbol{w}} \rangle \to m^{\star}, \ \frac{1}{d} \| \hat{\boldsymbol{w}} \|_{2}^{2} \to q^{\star}, \ \frac{1}{d} \| \hat{\boldsymbol{w}} \|_{p^{\star}}^{p^{\star}} \to P^{\star},$$

Remark 3.6. An immediate observation from the above equations is that E_{gen} is monotonically *increasing* as a function of the cosine of the angle between teacher and student $(m^*/\sqrt{\rho q^*})$, while E_{bnd} is *decreasing*. This has been observed before for boundary based classifiers (Tanay & Griffin, 2016; Tanner et al., 2024).

Theorem 3.5 therefore states that in order to characterize the robust generalization error in the high-dimensional limit, it is enough to compute three low-dimensional statistics of the RERM solution. Our next result shows that these quantities can be asymptotically computed without having to actually solve the high-dimensional minimization problem in eq. (6). **Theorem 3.7** (Self-consistent equations for ℓ_p norms). Under the same assumptions as Theorem 3.5, the summary statistics (m^*, q^*, P^*) are the unique solution of the following set of self-consistent equations:

$$\begin{cases} \hat{m} = \alpha \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} \mathrm{d}y \, \partial_{\omega} \mathcal{Z}_{0} f_{g}(\sqrt{q}\xi, P) \right] \\ \hat{q} = \alpha \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} \mathrm{d}y \, \mathcal{Z}_{0} f_{g}^{2}(\sqrt{q}\xi, P) \right] \\ \hat{V} = -\alpha \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} \mathrm{d}y \, \mathcal{Z}_{0} \partial_{\omega} f_{g}(\sqrt{q}\xi, P) \right] \\ \hat{P} = \varepsilon \alpha p^{\star} P^{-\frac{1}{p}} \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} \mathrm{d}y \, \mathcal{Z}_{0} y f_{g}(\sqrt{q}\xi, P) \right] \end{cases}$$

$$(8)$$

and

$$\begin{cases} m = \mathbb{E}_{\xi} \left[\partial_{\gamma} \mathcal{Z}_{w} f_{w}(\sqrt{q}\xi, \hat{P}, \lambda) \right] \\ q = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w} f_{w}(\sqrt{q}\xi, \hat{P}, \lambda)^{2} \right] \\ V = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w} \partial_{\gamma} f_{w}(\sqrt{q}\xi, \hat{P}, \lambda) \right] \\ P = \mathbb{E}_{\xi} \left[\mathcal{Z}_{w} \partial_{\hat{P}} \mathcal{M}_{\frac{\lambda}{\hat{V}} |\cdot|^{r} + \frac{\hat{P}}{\hat{V}} |\cdot|^{p^{\star}}} \left(\frac{\sqrt{\hat{q}}\xi}{\hat{V}} \right) \right] \end{cases}$$
(9)

where $\mathcal{Z}_{w} = \mathcal{Z}_{w}(\hat{m}\xi/\sqrt{\hat{q}}, \hat{m}/\sqrt{\hat{q}}), \ \mathcal{Z}_{0} = \mathcal{Z}_{0}(y, m\xi/\sqrt{q}, \rho - m^{2}/q)$ and $\xi \sim \mathcal{N}(0, 1)$, and:

$$\mathcal{Z}_0(y,\omega,V) = \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left[P_{\text{out}}(y \mid \sqrt{V}z + \omega) \right], \quad (10)$$

$$f_g(\omega, \hat{P}) = \left(\mathcal{P}_{Vg(y, \cdot -y\varepsilon \ {^p}\sqrt{P})}(\omega) - \omega\right)/V, \qquad (11)$$

$$\mathcal{Z}_{\mathbf{w}}(\gamma, \Lambda) = \mathbb{E}_{w \sim P_w} \left[e^{-\frac{1}{2}\Lambda w^2 + \gamma w} \right], \tag{12}$$

$$f_w(\gamma, \hat{P}, \Lambda) = \mathcal{P}_{\frac{\lambda}{\Lambda}|\cdot|^r + \frac{\hat{P}}{\Lambda}|\cdot|^{p^\star}}\left(\frac{\gamma}{\Lambda}\right).$$
(13)

where we indicate the proximal of a function $f : \mathbb{R} \to \mathbb{R}$ as $\mathcal{P}_{Vf(\cdot)}(\omega)$ and its moreau envelope with $\mathcal{M}_{Vf(\cdot)}(\omega)$.

Two remarks on these two results are in order.

Remark 3.8. Both results hold for any separable convex regularizer in the definition of the empirical risk in eq. (6). This is in contrast to many prior works in this field, which primarily consider ℓ_2 regularizations.

Remark 3.9. The first four equations (eq. (8)) depend only on the noise distribution and the loss function, while the second set (eq. (9)) depends on the regularization function and the dual norm of the perturbation.

3.2. Results for Mahalanobis norms

5 While the ℓ_p norm is the most frequently discussed in the robust learning literature, ℓ_p perturbations are isotropic, treating all covariates equally. Under the isotropic Gaussian Assumption 3.4, this is justified. However, it can be limiting under more realistic scenarios where the covariates are *structured*, and for instance some features are more relevant than others. Recently, (Tanner et al., 2024) introduced a model for studying adversarial training under structured convariates which considers perturbations under a *Mahalanobis* norm, allowing to weight the perturbation along different directions. However, the discussion in that work focused only on ℓ_2 regularization.

Since our goal in this work is to study what is the best regularization choice for a given perturbation geometry, we now derive asymptotic results akin to the ones of Section 3.1 under any combination of Mahalanobis perturbation and regularization norm. As before, we start by introducing our assumptions.

Assumption 3.10 (Mahalanobis norms). Given a positive definite matrix Σ_{δ} , we consider perturbations under a Mahalanobis norm $||x||_{\Sigma_{\delta}} = \sqrt{x^{\top}\Sigma_{\delta}x}$. Additionally, we consider the regularization function to be $r(w) = 1/2 w^{\top}\Sigma_{w}w$ for a positive definite matrix Σ_{w} .

Assumption 3.11 (Structured data). For each $i \in [n]$, the covariates $x_i \in \mathbb{R}^d$ are drawn i.i.d. from the data distribution $P_{in}(x) = \mathcal{N}_x(0, \Sigma_x)$. Then the corresponding y_i is sampled independently from the conditional distribution \mathbb{P}_{out} defined in eq. (1). The target weight vector $w_{\star} \in \mathbb{R}^{d}$ is drawn from a prior probability distribution $w_{\star} \sim P_{w} = \mathcal{N}_{w_{\star}}(\mathbf{0}, \boldsymbol{\Sigma}_{\theta})$, which we assume has limiting Mahalanobis norm given by $\rho = \lim_{d\to\infty} \mathbb{E}[\frac{1}{d}\boldsymbol{w}_{\star}^{\top}\boldsymbol{\Sigma}_{\boldsymbol{x}}\boldsymbol{w}_{\star}]$. Assumption 3.12 (Scaling of Adversarial Norm Constraint). The value of ε scales with the dimension d such that $\varepsilon/\sqrt{d} = O(1)$.

Assumption 3.13 (Convergence of spectra). We suppose that $\Sigma_x, \Sigma_\delta, \Sigma_\theta, \Sigma_w$ are simultaneously diagonalisable. We call $\Sigma_x = S^\top \operatorname{diag}(\omega_i)S$, $\Sigma_\delta = S^\top \operatorname{diag}(\zeta_i)S$ and $\Sigma_w = S^\top \operatorname{diag}(w_i)S$. We define $\overline{\theta} = S\Sigma_x^\top w_* / \sqrt{\rho}$. We assume that the empirical distributions of eigenvalues and the entries of $\overline{\theta}$ jointly converge to a probability distribution μ as

$$\sum_{i=1}^{d} \delta(\bar{\boldsymbol{\theta}}_{i} - \bar{\boldsymbol{\theta}}) \delta(\omega_{i} - \omega) \delta(\zeta_{i} - \zeta) \delta(w_{i} - w) \to \mu.$$
(14)

As in Section 3.1, our first result concerns the limiting robust error.

Theorem 3.14 (Limiting errors for Mahalanobis norm). Let $\hat{w}(S) \in \mathbb{R}^d$ denote the unique solution of the RERM problem in eq. (6). Then, under Assumptions 3.1 and 3.10 to 3.13, the standard, robust and boundary generalization error of \hat{w} converge to the following deterministic quantities

$$E_{\text{gen}}(\hat{\boldsymbol{w}}) = \frac{1}{\pi} \arccos\left(m^* / \sqrt{(\rho + \tau^2)q^*}\right),$$

$$E_{\text{bnd}}(\hat{\boldsymbol{w}}) = \int_0^{\varepsilon_g \frac{\sqrt{\rho^*}}{\sqrt{q^*}}} \operatorname{erfc}\left(\frac{-\frac{m^*}{\sqrt{q^*}}\nu}{\sqrt{2(\rho + \tau^2 - m^{*2}/q^*)}}\right) \frac{e^{-\frac{\nu^2}{2}}}{\sqrt{2\pi}} \,\mathrm{d}\nu,$$

$$E_{\text{rob}}(\hat{\boldsymbol{w}}) = E_{\text{gen}}(\hat{\boldsymbol{w}}) + E_{\text{bnd}}(\hat{\boldsymbol{w}})$$

where m^*, q^*, P^* are the limiting values of the following summary statistics:

$$\frac{\boldsymbol{w}_{\star}^{\top}\boldsymbol{\Sigma}_{\boldsymbol{x}}\hat{\boldsymbol{w}}}{d} \to \boldsymbol{m}^{\star}, \ \frac{\hat{\boldsymbol{w}}^{\top}\boldsymbol{\Sigma}_{\boldsymbol{x}}\hat{\boldsymbol{w}}}{d} \to \boldsymbol{q}^{\star}, \ \frac{\hat{\boldsymbol{w}}^{\top}\boldsymbol{\Sigma}_{\boldsymbol{\delta}}\hat{\boldsymbol{w}}}{d} \to \boldsymbol{P}^{\star}$$

As in Section 3.1, our next result shows that the summary statistics characterizing the limiting errors can be obtained from of a set of self-consistent equations.

Theorem 3.15 (Self-Consistent equations for Mahalanobis norm). Under the same assumptions as Theorem 3.14, the summary statistics (m^*, q^*, P^*) are the unique solution of the following set of self-consistent equations:

$$\begin{cases} \hat{m} = \alpha \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} dy \, \partial_{\omega} \mathcal{Z}_{0} f_{g}(\sqrt{q}\xi, P) \right] \\ \hat{q} = \alpha \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} dy \, \mathcal{Z}_{0} f_{g}^{2}(\sqrt{q}\xi, P) \right] \\ \hat{V} = -\alpha \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} dy \, \mathcal{Z}_{0} \partial_{\omega} f_{g}(\sqrt{q}\xi, P) \right] \\ \hat{P} = 2\varepsilon \alpha P^{-\frac{1}{2}} \mathbb{E}_{\xi} \left[\int_{\mathbb{R}} dy \, \mathcal{Z}_{0} y f_{g}(\sqrt{q}\xi, P) \right] \end{cases}$$

$$(15)$$

and

$$\begin{cases} m = \mathbb{E}_{\mu} \left[\frac{\hat{m} \theta^{2}}{\lambda w + \hat{V} \omega + \hat{P} \delta} \right] \\ q = \mathbb{E}_{\mu} \left[\frac{\hat{m}^{2} \hat{\theta}^{2} \omega + \hat{q} \omega^{2}}{(\lambda w + \hat{V} \omega + \hat{P} \delta)^{2}} \right] \\ V = \mathbb{E}_{\mu} \left[\frac{\omega}{\lambda w + \hat{V} \omega + \hat{P} \delta} \right] \\ P = \mathbb{E}_{\mu} \left[\zeta \frac{\hat{m}^{2} \hat{\theta}^{2} + \hat{q} \omega^{2}}{(\lambda w + \hat{V} \omega + \hat{P} \delta)^{2}} \right] \end{cases}$$
(16)

where μ is the joint limiting distribution for the spectrum of all the matrices.

Remark 3.16. Notice that the first set of equations is the same as in Theorem 3.7, as they depend only on the marginal distribution \mathbb{P}_{out} and the loss function.

Remark 3.17 (Interpretation of the self-consistent equations).The self-consistent equations in Theorems 3.7 and 3.15reduce the high-dimensional optimization problem in eq. (6)to tracking only three scalar statistics (m^*, q^*, P^*) . Theseequations, derived via Gordon's Min-Max theorem, fullycharacterize the asymptotic behavior of regularized RERMthrough the interaction of these quantities at optimality.

288 While the self-consistent equations in Theorem 3.7 and The-289 orem 3.15 do not admit a closed-form solution, they can be 290 efficiently solved using an iteration scheme (Appendix D). 291 Solving them yields precise curves for the generalization 292 errors of the final predictor as a function of the sample com-293 plexity α and regularization geometry, allowing us to draw 294 conclusions for the interplay between the regularization and 295 perturbation – see simulations in Section 5. 296

The details of the proofs of Theorems 3.5, 3.7, 3.14 and 3.15 are discussed in Appendix A. They are based on an adaptation of Gordon's Min-Max Theorem for convex empirical risk minimization problems (Thrampoulidis et al., 2014; Loureiro et al., 2021).

4. Which Regularization to Choose?

302

324

325

327

328

329

Our results in the previous section provide tight predictions 305 on the robust and standard generalization error of the set 306 of minimizers of the robust (regularized) empirical risk. 307 However, since the self-consistent equations describing the 308 robust errors are not closed, it is not straightforward to read 309 why some regularizers might produce better results than 310 others. In this section, we derive complementary uniform 311 convergence bounds based on the *Rademacher Complexity* 312 for linear predictors under various geometries. While these 313 bounds might not be numerically tight, they are distribution-314 agnostic, and provide a-priori guarantees for the error of a 315 predictor which are qualitatively useful. We start by intro-316 ducing concepts in a general way, before deriving guarantees 317 for the case considered in Section 3.2. 318

Let $\mathcal{H}_{\tilde{r}}$ be a hypothesis class of linear predictors of restricted complexity, as captured by a function $\tilde{r} : \mathbb{R}^d \to \mathbb{R}$. This function \tilde{r} plays the role of a regularizer, as in Section 3. We define:

$$\mathcal{H}_{\widetilde{r}} = \{ \mathbf{x} \to \langle \boldsymbol{w}, \mathbf{x} \rangle : \widetilde{r}(\boldsymbol{w}) \le \mathcal{W}_{\widetilde{r}}^2 \},$$
(17)

where $W_{\tilde{r}} > 0$ is an arbitrary upper bound. Central to the analysis of the generalization error uniformly inside the hypothesis class $\mathcal{H}_{\tilde{r}}$ is the notion of the (empirical) Rademacher Complexity (Koltchinskii, 2001) of $\mathcal{H}_{\tilde{r}}$:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\mathcal{H}_{\widetilde{r}}) = \mathbb{E}_{\sigma} \left[\frac{1}{n} \sup_{\boldsymbol{w}: \widetilde{r}(\boldsymbol{w}) \leq \mathcal{W}_{\widetilde{r}}^{2}} \sum_{i=1}^{n} \sigma_{i} \left\langle \boldsymbol{w}, \boldsymbol{x}_{i} \right\rangle \right], \quad (18)$$

where the σ_i 's are either -1 or 1 with equal probability. In the case of robust generalization with respect to $\|\cdot\|$ limited perturbations, it suffices to analyse the *worst-case* Rademacher Complexity of $\mathcal{H}_{\vec{r}}$:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\widetilde{r}}) = \mathbb{E}_{\sigma} \left[\frac{1}{n} \sup_{\boldsymbol{w}: \widetilde{r}(\boldsymbol{w}) \leq \mathcal{W}_{\widetilde{r}}^2} \sum_{i=1}^n \sigma_i \min_{\|\boldsymbol{\delta}_i\| \leq \varepsilon} \langle \boldsymbol{w}, \boldsymbol{x}_i + \boldsymbol{\delta}_i \rangle \right]$$

With these ingredients in place, we can state the following bound on the robust generalization gap of any predictor in $\mathcal{H}_{\tilde{r}}$.

Theorem 4.1 (Mohri et al. (2012); Awasthi et al. (2020)). For any $\delta > 0$, with probability at least $1 - \delta$ over the draw of the dataset S, for all $w \in \mathbb{R}^d$ such that $r(w) \leq W_r^2$ (eq. (17)), it holds that

$$E_{\rm rob}(\boldsymbol{w}) \leq \hat{E}_{\rm rob}(\boldsymbol{w}) + 2\,\hat{\mathfrak{R}}_{\rm S}(\widetilde{\mathcal{H}}_r) + 3\sqrt{\frac{\log 2/\delta}{2n}},$$
 (19)

where

$$\hat{E}_{\rm rob}(\boldsymbol{w}) = \frac{1}{n} \sum_{i=1}^{n} \max_{\|\boldsymbol{\delta}_i\| \le \varepsilon} \mathbb{1}(y_i \hat{y}(\boldsymbol{w}, \boldsymbol{x}_i + \boldsymbol{\delta}_i) \le 1) \quad (20)$$

is a robust empirical error.

Theorem 4.1 promises that a tight bound on the worst-case Rademacher complexity of $\mathcal{H}_{\tilde{r}}$ can bound the (robust) generalization gap of any predictor in $\mathcal{H}_{\tilde{r}}$. The next Proposition realises this goal for the general class of *strongly convex* functions \tilde{r} . This will permit the study of the cases of Section 3.2.

Proposition 4.2. Let $\varepsilon, \sigma > 0$. Consider a sample $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$, and let $\mathcal{H}_{\tilde{r}}$ be the hypothesis class defined in eq. (17), where \tilde{r} is σ -strongly convex with respect to a norm $_{r} \|\cdot\|$. Then, it holds:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\widetilde{r}}) \leq \max_{i \in [n]} {}_{r} \|\mathbf{x}_{i}\|_{\star} \mathcal{W}_{\widetilde{r}} \sqrt{\frac{2}{\sigma n}} + \frac{\varepsilon}{2\sqrt{n}} \sup_{\boldsymbol{w}: \widetilde{r}(\boldsymbol{w}) \leq \mathcal{W}_{\widetilde{r}}^{2}} \|\boldsymbol{w}\|_{\star},$$
(21)

where $_{r}\|\cdot\|_{\star}, \|\cdot\|_{\star}$ denote the dual norms of $_{r}\|\cdot\|, \|\cdot\|$, respectively.

The proof (Appendix B) leverages a fundamental result on (standard) Rademacher complexities for strongly convex functions due to Kakade et al. (2008) and a symmetrization argument.

This result informs us that the worst-case Rademacher complexity can decompose into two terms: one which characterizes the standard error and one that scales with the magnitude of perturbation ε and depends on the *dual* norm of the perturbation. Thus, we expect that a regularization which promotes a small second term in the RHS of eq. (21) will likely mean a smaller robust generalization gap, as ε increases. This can be further elucidated in the following subcases (proofs appear in appendix B), for which we already derived exact asymptotics in Section 3.2:

||·|| = ||·||_{Σδ} and *ĩ*(*w*) = ||*w*||₂²: this corresponds to perturbations with respect to a a symmetric positive definite matrix Σ_δ ∈ ℝ^{d×d}, while we regularize in the Euclidean norm. In this case, we obtain:

337

338

339 340

341

342

343

345

346

347

349

350

351

352

353

354

355

356

358

359

360

361

362

363

366

367

368

369

370

371

373

374

375

376

377

Corollary 4.3. Let $\varepsilon > 0$ and symmetric positive definite $\Sigma_{\delta} \in \mathbb{R}^{d \times d}$. Then:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\|\cdot\|_{2}^{2}}) \leq \frac{\max_{i \in [n]} \|\mathbf{x}_{i}\|_{2} \mathcal{W}_{2}}{\sqrt{n}} + \frac{\varepsilon \mathcal{W}_{2}}{2\sqrt{n}} \sqrt{\lambda_{\min}^{-1}(\boldsymbol{\Sigma}_{\boldsymbol{\delta}})}.$$

||·|| = ||·||_{Σδ} and *r̃*(*w*) = ||*w*||²_{Σw}: this corresponds to perturbations with respect to a symmetric positive definite matrix Σ_δ ∈ ℝ^{d×d} and regularization with respect to a norm induced by another matrix Σ_w ∈ ℝ^{d×d}. We will analyze the special case where Σ_δ and Σ_w share the same set of eigenvectors.

Corollary 4.4. Let $\varepsilon > 0$. Let $\Sigma_{\boldsymbol{w}} = \sum_{i=1}^{d} \alpha_i \mathbf{v}_i \mathbf{v}_i^T$ and $\Sigma_{\boldsymbol{\delta}} = \sum_{i=1}^{d} \lambda_i \mathbf{v}_i \mathbf{v}_i^T$, with $\mathbf{v}_i \in \mathbb{R}^d$ being orthonormal. Then:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\|\cdot\|_{A}^{2}}) \leq \frac{\mathcal{W}_{A} \max_{i \in [n]} \|\mathbf{x}_{i}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}^{-1}}}{\sqrt{n}} + \frac{\varepsilon \mathcal{W}_{A}}{2\sqrt{n}} \sqrt{\max_{i \in [d]} \frac{1}{\lambda_{i} \alpha_{i}}}.$$
(22)

Hence, we deduce that regularizing the class of linear predictors with $\Sigma_w = \Sigma_{\delta}^{-1}$, where Σ_{δ} is the matrix of the perturbation norm, can more effectively control the robust generalization error.

Similar results can be derived in the context of ℓ_p perturbations - see Yin et al. (2019); Awasthi et al. (2020) and Appendix B. In fact, mirroring our analysis, the robust generalization error there is controlled by the $\|\cdot\|_{p^*}$ norm of the weights. We explore the effect of the regularizer numerically with simulations next.

5. Experiments

Leveraging our exact results from Section 3 and guided by
the predictions of Section 4, in this Section we numerically
investigate the role of the regularization geometry in the
robustness and accuracy of robust empirical risk minimizers. Experimental details and further ablation studies are
in Appendix C.



Figure 1: Generalization error of RERMs in the low sample complexity regime under ℓ_{∞} perturbations for various choices of regularization. We see that the edge of ℓ_1 over the rest of the methods stems from the boundary error (E_{bnd}) which goes to zero as $\alpha \to 0^+$. Setting: $\varepsilon = 0.2$ and optimally tuned regularization parameter λ . The bullet points with the error bars are RERM simulations for d = 1000 (10 random seeds), while the solid lines correspond to the theoretical predictions.

5.1. Importance of Regularization in the Scarce Data Regime

First, we consider the setting of Section 3.1, with perturbations constrained in their ℓ_{∞} norms, for three different regularizers (ℓ_1 , ℓ_2 and ℓ_3 norms). Figure 1 compares the generalization errors of the solutions of eq. (6) for the various regularizers and plots them as a function of the sample complexity α . Note that when α is small (scarce data), the ℓ_1 regularized solution (dual norm of ℓ_{∞}) provides better defense against ℓ_{∞} perturbations. Interestingly, this is due to the fact that the boundary error approaches zero as $\alpha \rightarrow 0^+$, only in the case when r = 1 (Figure 3, bottom). We analytically explore this phenomenon further in Appendix A.8, where we analyze the boundary error from Theorem 3.5 and probe its dependence on the overlap parameters (m^*, q^*, P^*).

The phase diagram of Figure 2 further elucidates the difference of the methods as a function of ε . We display the difference in robust generalization error between ℓ_2 and ℓ_1 regularized solutions versus attack budget ε and sample complexity α , with *optimally tuned* regularization coefficient λ . We observe that ℓ_1 outperforms ℓ_2 regularization in regions of high ε and low α .

Figure 3 demonstrates $E_{\rm rob}$ in the structured case of Section 3.2, where the perturbations are constrained in a Mahalanobis norm $\|\cdot\|_{\Sigma_{\delta}}$. We observe that regularizing the weights of the solution with the dual norm of the perturbation $(\|\cdot\|_{\Sigma_{\delta}})$ yields better robustness, while the gap between the various methods increases as ε grows.



Figure 2: Difference between robust generalization errors for r = 2 and r = 1 as a function of ε and α for ℓ_{∞} attacks. Green zones correspond to areas where the the dual norm regularization is better than ℓ_2 .

399

400

401

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439



Figure 3: Difference between robust generalization error for Σ_{δ} perturbations. We see that a regularization with the dual norm has the best adversarial error for different choices of ε . The points with the error bars (std) are RERM simulations for d = 1000 (10 random seeds), while the solid lines correspond to the theoretical predictions.

5.2. Optimal Regularization Geometry as a Function of ε

While the previous figures compared the various regularizers r(w) as ε grows, it is not clear what exactly the relationship is between optimal r(w) and perturbation strength. In particular, we expect when $\varepsilon = 0$, ℓ_2 -regularized solutions to achieve better accuracy, due to the fact that the data are Gaussian. However, it is not clear how the transition to the dual norm happens.

We examine this relationship in Figure 4, where we plot the robust generalization error for various values of perturbation ε and regularization order r for a fixed value of sample complexity α . We observe that, as the attack strength increases, the order of the optimal regularization *smoothly* transitions



Figure 4: Robust generalization error of the solution of regularized RERM as a function of the regularization order r, i.e. $r(w) = \lambda ||w||_r^r$ for various perturbations strengths ε . Sample complexity $\alpha = 1.0$. Regularization coefficients λ are optimally tuned. The inside figure shows how the optimal value of r scales with ε .

from r = 2 to r = 1. Hence, there is a regime of perturbation scale ε where neither r = 2 nor r = 1 is optimal, but an order of $r \in (1, 2)$ achieves the least robust test error.

6. Conclusion

We studied the role of regularization in robust empirical risk minimization (adversarial training) for a variety of perturbation and regularization norms. We derived an exact asymptotic description of the robust and standard generalization error in the high-dimensional proportional limit, and we showed results for the (worst-case) Rademacher Complexity of linear predictors in the case of structured perturbations. Phase diagrams and exact scaling laws, afforded by our analysis, suggest that choosing the right regularization becomes increasingly important as ε grows, and, in fact, this optimal regularization often corresponds to the dual norm of the perturbation. Furthermore, our results reveal a curious, smooth, transition between different optimal regularizations (ℓ_2 to ℓ_1) with increasing perturbation strength; a phenomenon that has not yet been captured by any other theoretical work.

It would be interesting for future work to investigate the interplay between regularization and perturbation geometry in non-linear models, such as the random features model (Mei & Montanari, 2022; Gerace et al., 2021; Hassani & Javanmard, 2024).

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

440 **References**

474

475

476

477

484

- Adomaityte, U., Defilippis, L., Loureiro, B., and Sicuro,
 G. High-dimensional robust regression under heavytailed data: Asymptotics and universality. *arXiv preprint arXiv:2309.16476*, 2023.
- Adomaityte, U., Sicuro, G., and Vivo, P. Classification of heavy-tailed features in high dimensions: a superstatistical approach. *Advances in Neural Information Processing Systems*, 36, 2024.
- Aubin, B., Krzakala, F., Lu, Y., and Zdeborová, L. Generalization error in high-dimensional perceptrons: Approaching bayes error with convex optimization. In Larochelle,
 H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H.
 (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 12199–12210. Curran Associates,
 Inc., 2020.
- Awasthi, P., Frank, N., and Mohri, M. Adversarial learning guarantees for linear hypotheses and neural networks. In *International Conference on Machine Learning*, pp. 431–441. PMLR, 2020.
- 463 Bean, D., Bickel, P. J., El Karoui, N., and Yu, B. Optimal mestimation in high-dimensional regression. *Proceedings of the National Academy of Sciences*, 110(36):14563– 14568, 2013a.
- Bean, D., Bickel, P. J., Karoui, N. E., and Yu, B.
 Optimal m-estimation in high-dimensional regression. *Proceedings of the National Academy of Sciences*,
 110(36):14563–14568, 2013b. doi: 10.1073/pnas.
 1307845110. URL https://www.pnas.org/doi/
 abs/10.1073/pnas.1307845110.
 - Bhagoji, A. N., Cullina, D., and Mittal, P. Lower bounds on adversarial robustness from optimal transport. *Advances in Neural Information Processing Systems*, 32, 2019.
- Bordelon, B., Canatar, A., and Pehlevan, C. Spectrum dependent learning curves in kernel regression and wide neural networks. In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1024–1034. PMLR, 13–18 Jul 2020.
- Boyd, S. and Vandenberghe, L. *Convex Optimization*.
 Cambridge University Press, 2004. doi: 10.1017/
 CBO9780511804441.
- Carmon, Y., Raghunathan, A., Schmidt, L., Duchi, J. C., and Liang, P. Unlabeled data improves adversarial robustness. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing*

Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 11190–11201, 2019.

- Chen, T., Zhang, Z., Liu, S., Chang, S., and Wang, Z. Robust overfitting may be mitigated by properly learned smoothening. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021, 2021.
- Cortes, C. and Vapnik, V. Support-vector networks. *Ma-chine Learning*, 20(3):273–297, 1995.
- Dan, C., Wei, Y., and Ravikumar, P. Sharp statistical guaratees for adversarially robust gaussian classification. In *International Conference on Machine Learning*, pp. 2345– 2355. PMLR, 2020.
- Dandi, Y., Stephan, L., Krzakala, F., Loureiro, B., and Zdeborová, L. Universality laws for gaussian mixtures in generalized linear models. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 54754–54768. Curran Associates, Inc., 2023.
- Dohmatob, E. and Scetbon, M. Precise accuracy/robustness tradeoffs in regression: Case of general norms. In *Forty-first International Conference on Machine Learning*, 2024.
- Donoho, D. and Tanner, J. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 367(1906): 4273–4293, 2009.
- Dudeja, R., M. Lu, Y., and Sen, S. Universality of approximate message passing with semirandom matrices. *The Annals of Probability*, 51(5):1616–1683, 2023.
- Gerace, F., Loureiro, B., Krzakala, F., Mézard, M., and Zdeborová, L. Generalisation error in learning with random features and the hidden manifold model*. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124013, dec 2021. doi: 10.1088/1742-5468/ ac3ae6. URL https://dx.doi.org/10.1088/ 1742-5468/ac3ae6.
- Gerace, F., Krzakala, F., Loureiro, B., Stephan, L., and Zdeborová, L. Gaussian universality of perceptrons with random labels. *Phys. Rev. E*, 109: 034305, Mar 2024. doi: 10.1103/PhysRevE.109. 034305. URL https://link.aps.org/doi/10. 1103/PhysRevE.109.034305.

- Goldt, S., Loureiro, B., Reeves, G., Krzakala, F., Mezard,
 M., and Zdeborová, L. The gaussian equivalence of generative models for learning with shallow neural networks.
 Proceedings of Machine Learning Research. 145, pp. 426–
 471, 2021.
- 500 Goldt, S., Loureiro, B., Reeves, G., Krzakala, F., Mezard, 501 M., and Zdeborova, L. The gaussian equivalence of gen-502 erative models for learning with shallow neural networks. 503 In Bruna, J., Hesthaven, J., and Zdeborova, L. (eds.), Pro-504 ceedings of the 2nd Mathematical and Scientific Machine 505 Learning Conference, volume 145 of Proceedings of Ma-506 chine Learning Research, pp. 426-471. PMLR, 16-19 507 Aug 2022. 508
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. In Bengio, Y. and LeCun, Y. (eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- Gordon, Y. On milman's inequality and random subspaces which escape through a mesh in Rⁿ. In Lindenstrauss, J. and Milman, V. D. (eds.), *Geometric Aspects of Functional Analysis*, pp. 84–106, Berlin, Heidelberg, 1988.
 Springer Berlin Heidelberg. ISBN 978-3-540-39235-4.
- Gunasekar, S., Lee, J. D., Soudry, D., and Srebro, N. Characterizing implicit bias in terms of optimization geometry. In Dy, J. G. and Krause, A. (eds.), *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, pp. 1827–1836. PMLR, 2018.
- Hassani, H. and Javanmard, A. The curse of overparametrization in adversarial training: Precise analysis
 of robust generalization for random features regression. *The Annals of Statistics*, 52(2):441–465, 2024.
- Hu, H. and Lu, Y. M. Universality laws for high-dimensional learning with random features. *IEEE Transactions on Information Theory*, 69(3):1932–1964, 2023. doi: 10. 1109/TIT.2022.3217698.
- Isserlis, L. On a formula for the product-moment coefficient
 of any order of a normal frequency distribution in any
 number of variables. *Biometrika*, 12(1/2):134–139, 1918.

- Javanmard, A. and Soltanolkotabi, M. Precise statistical analysis of classification accuracies for adversarial training. *The Annals of Statistics*, 50(4):2127–2156, 2022.
- Kakade, S. M., Sridharan, K., and Tewari, A. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In Koller, D., Schuurmans, D., Bengio,

Y., and Bottou, L. (eds.), Advances in Neural Information Processing Systems 21, Proceedings of the Twenty-Second Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 8-11, 2008, pp. 793–800. Curran Associates, Inc., 2008.

- Koltchinskii, V. Rademacher penalties and structural risk minimization. *IEEE Trans. Inf. Theory*, 47(5):1902–1914, 2001.
- Krogh, A. and Hertz, J. A simple weight decay can improve generalization. Advances in neural information processing systems, 4, 1991.
- Liao, Z., Couillet, R., and Mahoney, M. W. A random matrix analysis of random fourier features: beyond the gaussian kernel, a precise phase transition, and the corresponding double descent. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 13939–13950. Curran Associates, Inc., 2020. URL https://proceedings.neurips. cc/paper_files/paper/2020/file/ a03fa30821986dff10fc66647c84c9c3-Paper. pdf.
- Loureiro, B., Gerbelot, C., Cui, H., Goldt, S., Krzakala, F., Mezard, M., and Zdeborová, L. Learning curves of generic features maps for realistic datasets with a teacherstudent model. *Advances in Neural Information Processing Systems*, 34:18137–18151, 2021.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*, 2018.
- Mai, X., Liao, Z., and Couillet, R. A large scale analysis of logistic regression: Asymptotic performance and new insights. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing* (*ICASSP*), pp. 3357–3361, 2019. doi: 10.1109/ICASSP. 2019.8683376.
- Mei, S. and Montanari, A. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 75(4):667–766, 2022.
- Mignacco, F., Krzakala, F., Lu, Y., Urbani, P., and Zdeborova, L. The role of regularization in classification of high-dimensional noisy gaussian mixture. In *International conference on machine learning*, pp. 6874–6883. PMLR, 2020.

Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Founda- tions of Machine Learning*. Adaptive computation and
machine learning. MIT Press, 2012. ISBN 978-0-26201825-8.

554

562

588

589

590

591

592

- Montanari, A. and Saeed, B. N. Universality of empirical risk minimization. In Loh, P.-L. and Raginsky,
 M. (eds.), *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pp. 4310–4312. PMLR, 02–05 Jul
 2022. URL https://proceedings.mlr.press/ v178/montanari22a.html.
- Okajima, K., Meng, X., Takahashi, T., and Kabashima, 563 Y. Average case analysis of lasso under ultra 564 sparse conditions. In Ruiz, F., Dy, J., and van de 565 Meent, J.-W. (eds.), Proceedings of The 26th Interna-566 tional Conference on Artificial Intelligence and Statis-567 tics, volume 206 of Proceedings of Machine Learn-568 ing Research, pp. 11317-11330. PMLR, 25-27 Apr 569 2023. URL https://proceedings.mlr.press/ 570 v206/okajima23a.html. 571
- 572
 573 Parikh, N. and Boyd, S. Proximal algorithms. *Found. Trends*574 *Optim.*, 1(3):127–239, jan 2014. ISSN 2167-3888. doi:
 575 10.1561/2400000003. URL https://doi.org/10.
 576 1561/240000003.
- 577 Pesce, L., Krzakala, F., Loureiro, B., and Stephan, L. Are 578 Gaussian data all you need? The extents and limits of 579 universality in high-dimensional generalized linear es-580 timation. In Krause, A., Brunskill, E., Cho, K., En-581 gelhardt, B., Sabato, S., and Scarlett, J. (eds.), Pro-582 ceedings of the 40th International Conference on Ma-583 chine Learning, volume 202 of Proceedings of Machine 584 Learning Research, pp. 27680-27708. PMLR, 23-29 Jul 585 2023. URL https://proceedings.mlr.press/ 586 v202/pesce23a.html. 587
 - Raghunathan, A., Xie, S. M., Yang, F., Duchi, J., and Liang, P. Understanding and mitigating the tradeoff between robustness and accuracy. *arXiv preprint arXiv:2002.10716*, 2020.
- Rebuffi, S., Gowal, S., Calian, D. A., Stimberg, F., Wiles,
 O., and Mann, T. A. Fixing data augmentation to improve adversarial robustness. *CoRR*, abs/2103.01946, 2021.
 URL https://arxiv.org/abs/2103.01946.
- Rice, L., Wong, E., and Kolter, J. Z. Overfitting in adversarially robust deep learning. 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*, pp. 8093–8104. PMLR, 2020.

- Schröder, D., Cui, H., Dmitriev, D., and Loureiro, B. Deterministic equivalent and error universality of deep random features learning. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 30285–30320. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/schroder23a.html.
- Seung, H. S., Sompolinsky, H., and Tishby, N. Statistical mechanics of learning from examples. *Phys. Rev. A*, 45: 6056–6091, Apr 1992.
- Shafahi, A., Najibi, M., Ghiasi, A., Xu, Z., Dickerson, J. P., Studer, C., Davis, L. S., Taylor, G., and Goldstein, T. Adversarial training for free! In Wallach, H. M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 3353–3364, 2019.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient descent on separable data. J. Mach. Learn. Res., 19:70:1–70:57, 2018.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. J., and Fergus, R. Intriguing properties of neural networks. In Bengio, Y. and LeCun, Y. (eds.), 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014. URL http://arxiv.org/abs/1312.6199.
- Taheri, H., Pedarsani, R., and Thrampoulidis, C. Asymptotic behavior of adversarial training in binary linear classification. *IEEE Trans. Neural Netw. Learn. Syst.*, PP, July 2023.
- Tanay, T. and Griffin, L. A boundary tilting persepective on the phenomenon of adversarial examples. *arXiv preprint arXiv:1608.07690*, 2016.
- Tanner, K., Vilucchio, M., Loureiro, B., and Krzakala, F. A high dimensional model for adversarial training: Geometry and trade-offs. arXiv preprint arXiv:2402.05674, 2024.
- Tao, T. and Vu, V. Random matrices: Universality of local eigenvalue statistics up to the edge. *Communications in Mathematical Physics*, 298:549–572, 2010.
- Thrampoulidis, C., Oymak, S., and Hassibi, B. The gaussian min-max theorem in the presence of convexity. *arXiv preprint arXiv:1408.4837*, 2014.

605 606 607 608 609 610	Thrampoulidis, C., Oymak, S., and Hassibi, B. Regularized Linear Regression: A precise analysis of the estimation error. In Grünwald, P., Hazan, E., and Kale, S. (eds.), <i>Proceedings of The 28th Conference on Learning Theory</i> , volume 40 of <i>Proceedings of Machine Learning Research</i> , pp. 1683–1709, Paris, France, 03–06 Jul 2015. PMLR.	Honolulu, Ha Machine Lean 2023. Wei, A., Hu, W dom matrix r resentations s
 611 612 613 614 615 	Thrampoulidis, C., Abbasi, E., and Hassibi, B. Precise error analysis of regularized <i>m</i> -estimators in high dimensions. <i>IEEE Transactions on Information Theory</i> , 64(8):5592– 5628, 2018.	Song, L., Sze Proceedings of chine Learnin Learning Res.
616617618619	Tikhonov, A. On the solution of ill-posed problems and the method of regularization. <i>Doklady Akademii Nauk USSR</i> , pp. 501–504, 1963.	v162/wei2 Xiao, J., Sun, R
620 621 622 623	Tsilivis, N., Frank, N., Srebro, N., and Kempe, J. The price of implicit bias in adversarially robust generalization. <i>arXiv preprint arXiv:2406.04981</i> , 2024.	ization. In A Seventh Annu. - July 3, 2023
624 625 626 627 628	Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness may be at odds with accuracy. In 7th International Conference on Learning Representa- tions, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019, 2019.	Xiao, L., Hu, F ton, J. Precise for dot-produ
629 630 631 632	Vapnik, V. <i>Statistical learning theory</i> . Wiley, 1998. ISBN 978-0-471-03003-4.	Information F Yin, D., Ramch complexity for
 632 633 634 635 636 	Vapnik, V. N. and Chervonenkis, A. Y. On the uniform convergence of relative frequencies of events to their probabilities. <i>Theory of Probability & Its Applications</i> , 16(2):264–280, 1971.	Chaudhuri, K of the 36th In ing, ICML 20. USA, volume
 637 638 639 640 641 642 	Vilucchio, M., Troiani, E., Erba, V., and Krzakala, F. Asymptotic characterisation of the performance of robust linear regression in the presence of outliers. In <i>Interna-</i> <i>tional Conference on Artificial Intelligence and Statistics</i> , pp. 811–819. PMLR, 2024.	Zhai, R., Cai, T. Wang, L. Adv more unlabely 2019.
643 644 645 646 647	 Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, İ., 	Zhang, H., Yu, Jordan, M. 7 robustness an machine learn
648 649 650 651 652 653	Feng, Y., Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. <i>Nature Methods</i> , 17:261–272, 2020.	Zou, A., Wang, versal and trai guage models
654 655 656	Wang, Z., Pang, T., Du, C., Lin, M., Liu, W., and Yan, S. Better diffusion models further improve adversarial train- ing. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B.,	

Sabato, S., and Scarlett, J. (eds.), International Confer-

ence on Machine Learning, ICML 2023, 23-29 July 2023,

657

658

659

Honolulu, Hawaii, USA, volume 202 of *Proceedings of Machine Learning Research*, pp. 36246–36263. PMLR, 2023.

- Wei, A., Hu, W., and Steinhardt, J. More than a toy: Random matrix models predict how real-world neural representations generalize. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 23549–23588. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/ v162/wei22a.html.
- Xiao, J., Sun, R., Long, Q., and Su, W. Bridging the gap: Rademacher complexity in robust and standard generalization. In Agrawal, S. and Roth, A. (eds.), *The Thirty Seventh Annual Conference on Learning Theory, June 30* - July 3, 2023, Edmonton, Canada, volume 247 of Proceedings of Machine Learning Research, pp. 5074–5075. PMLR, 2024.
- Xiao, L., Hu, H., Misiakiewicz, T., Lu, Y., and Pennington, J. Precise learning curves and higher-order scalings for dot-product kernel regression. *Advances in Neural Information Processing Systems*, 35:4558–4570, 2022.
- Yin, D., Ramchandran, K., and Bartlett, P. L. Rademacher complexity for adversarially robust generalization. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *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, pp. 7085–7094. PMLR, 2019.
- Zhai, R., Cai, T., He, D., Dan, C., He, K., Hopcroft, J., and Wang, L. Adversarially robust generalization just requires more unlabeled data. arXiv preprint arXiv:1906.00555, 2019.
- Zhang, H., Yu, Y., Jiao, J., Xing, E., El Ghaoui, L., and Jordan, M. Theoretically principled trade-off between robustness and accuracy. In *International conference on machine learning*, pp. 7472–7482. PMLR, 2019.
- Zou, A., Wang, Z., Kolter, J. Z., and Fredrikson, M. Universal and transferable adversarial attacks on aligned language models. *CoRR*, abs/2307.15043, 2023.

nNumber of training samplesdInput dimension $\alpha = n/d$ Sample complexity (ratio of samples to dimension) $x_i \in \mathbb{R}^d$ Input features for sample i $y_i \in \{-1, +1\}$ Binary label for sample i $w_\star \in \mathbb{R}^d$ Teacher (true) weight vector $\hat{w} \in \mathbb{R}^d$ Learned weight vector (student) ε Adversarial perturbation budget $\ \cdot\ _p$ ℓ_p norm, defined as $\ x\ _p = (\sum_i x_i ^p)^{1/p}$ p^* Dual exponent of p , satisfying $1/p + 1/p^* = 1$ $E_{rob}(\hat{w})$ Robust generalization error $E_{gen}(\hat{w})$ Standard generalization error $E_{gen}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic dual norm of student weights P^* Asymptotic dual norm of student weights	Number of training samples Input dimension $= n/d$ Sample complexity (ratio of samples to dimension) $= n/d$ Sample complexity (ratio of samples to dimension) $i \in \mathbb{R}^d$ Input features for sample i $\in \{-1, +1\}$ Binary label for sample i $* \in \mathbb{R}^d$ Teacher (true) weight vector $\in \mathbb{R}^d$ Learned weight vector (student)Adversarial perturbation budget $ _p$ ℓ_p norm, defined as $ x _p = (\sum_i x_i ^p)^{1/p}$ $\operatorname{rob}(\hat{w})$ Robust generalization error $\operatorname{gen}(\hat{w})$ Standard generalization error $\operatorname{gen}(\hat{w})$ Boundary error (difference between robust and standard error) $\operatorname{Regularization strength parameter}$ w Regularization function \cdot)Surrogate loss function \cdot)Surrogate loss function w Prior distribution neacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization \star Asymptotic overlap between teacher and student weights \star Asymptotic dual norm of student weights \star Asymptotic dual norm of student weights \star Asymptotic dual norm of student weights \star Table 1: Notation Table	Symbol	Description
dInput dimension $\alpha = n/d$ Sample complexity (ratio of samples to dimension) $x_i \in \mathbb{R}^d$ Input features for sample i $y_i \in \{-1, +1\}$ Binary label for sample i $w_\star \in \mathbb{R}^d$ Teacher (true) weight vector $\hat{w} \in \mathbb{R}^d$ Learned weight vector (student) ε Adversarial perturbation budget $\ \cdot\ _p$ ℓ_p norm, defined as $\ x\ _p = (\sum_i x_i ^p)^{1/p}$ p^* Dual exponent of p , satisfying $1/p + 1/p^* = 1$ $E_{rob}(\hat{w})$ Robust generalization error $E_{gen}(\hat{w})$ Standard generalization error $E_{gen}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{uut} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights P^* Asymptotic dual norm of student weights P^* Asymptotic dual norm of student weights	Input dimension $= n/d$ Sample complexity (ratio of samples to dimension) $i \in \mathbb{R}^d$ Input features for sample i $\in \{-1, +1\}$ Binary label for sample i $* \in \mathbb{R}^d$ Teacher (true) weight vector $\in \mathbb{R}^d$ Learned weight vector (student)Adversarial perturbation budget $ _p$ ℓ_p norm, defined as $ x _p = (\sum_i x_i ^p)^{1/p}$ Dual exponent of p , satisfying $1/p + 1/p^* = 1$ $rob(\hat{w})$ Robust generalization error $gen(\hat{w})$ Standard generalization error $gen(\hat{w})$ Boundary error (difference between robust and standard error) $Regularization strength parameter$ w Regularization function \cdot)Surrogate loss function \cdot Output channel (conditional probability of labels) n Input distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization \star Asymptotic overlap between teacher and student weights \star Asymptotic dual norm of student weights	n	Number of training samples
$\alpha = n/d$ Sample complexity (ratio of samples to dimension) $x_i \in \mathbb{R}^d$ Input features for sample i $y_i \in \{-1, +1\}$ Binary label for sample i $w_\star \in \mathbb{R}^d$ Teacher (true) weight vector $\hat{w} \in \mathbb{R}^d$ Learned weight vector (student) ε Adversarial perturbation budget $\ \cdot\ _p$ ℓ_p norm, defined as $\ x\ _p = (\sum_i x_i ^p)^{1/p}$ p^* Dual exponent of p , satisfying $1/p + 1/p^* = 1$ $E_{rob}(\hat{w})$ Robust generalization error $E_{gen}(\hat{w})$ Standard generalization error $E_{bnd}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization strength parameter $\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weights P^* Asymptotic dual norm of student weights	$= n/d$ Sample complexity (ratio of samples to dimension) $g \in \mathbb{R}^d$ Input features for sample i $\in \{-1, +1\}$ Binary label for sample i $\star \in \mathbb{R}^d$ Teacher (true) weight vector $\in \mathbb{R}^d$ Learned weight vector (student)Adversarial perturbation budget $\ _p$ ℓ_p norm, defined as $\ x\ _p = (\sum_i x_i ^p)^{1/p}$ $\text{Trob}(\hat{w})$ Robust generalization error $\text{gen}(\hat{w})$ Standard generalization error $\text{gen}(\hat{w})$ Boundary error (difference between robust and standard error) $\text{Regularization strength parameter}$ w Regularization function.)Surrogate loss function out Output channel (conditional probability of labels) n Input distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Asymptotic overlap between teacher and student weights \star Asymptotic dual norm of student weights	d	Input dimension
$x_i \in \mathbb{R}^d$ Input features for sample i $y_i \in \{-1, +1\}$ Binary label for sample i $w_{\star} \in \mathbb{R}^d$ Teacher (true) weight vector $\hat{w} \in \mathbb{R}^d$ Learned weight vector (student) ε Adversarial perturbation budget $\ \cdot\ _p$ ℓ_p norm, defined as $\ x\ _p = (\sum_i x_i ^p)^{1/p}$ p^* Dual exponent of p , satisfying $1/p + 1/p^* = 1$ $E_{rob}(\hat{w})$ Robust generalization error $E_{gen}(\hat{w})$ Standard generalization error $E_{bnd}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization function $g(\cdot)$ Surrogate loss function \mathcal{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights p^* Asymptotic dual norm of student weights P^* Asymptotic dual norm of student weights	$\begin{array}{llllllllllllllllllllllllllllllllllll$	$\alpha = n/d$	Sample complexity (ratio of samples to dimension)
$y_i \in \{-1, +1\}$ Binary label for sample i $w_{\star} \in \mathbb{R}^d$ Teacher (true) weight vector $\hat{w} \in \mathbb{R}^d$ Learned weight vector (student) ε Adversarial perturbation budget $\ \cdot\ _p$ ℓ_p norm, defined as $\ x\ _p = (\sum_i x_i ^p)^{1/p}$ p^* Dual exponent of p , satisfying $1/p + 1/p^* = 1$ $E_{rob}(\hat{w})$ Robust generalization error $E_{gen}(\hat{w})$ Standard generalization error $E_{bnd}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization strength parameter $\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights P^* Asymptotic dual norm of student weights P^* Asymptotic dual norm of student weights	$ \{ \{-1, +1\} \} $ Binary label for sample <i>i</i> $ * \in \mathbb{R}^{d} $ Teacher (true) weight vector $ \in \mathbb{R}^{d} $ Learned weight vector (student) Adversarial perturbation budget $ * \ _{p} $ ℓ_{p} norm, defined as $\ x\ _{p} = (\sum_{i} x_{i} ^{p})^{1/p}$ Dual exponent of <i>p</i> , satisfying $1/p + 1/p^{*} = 1$ rob (\hat{w}) Robust generalization error gen (\hat{w}) Standard generalization error bind (\hat{w}) Boundary error (difference between robust and standard error) Regularization strength parameter w) Regularization function $\cdot)$ Surrogate loss function 0 Utput channel (conditional probability of labels) n Input distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization * Asymptotic overlap between teacher and student weights \star Asymptotic dual norm of student weights \star Asymptotic dual norm	$oldsymbol{x}_i \in \mathbb{R}^d$	Input features for sample <i>i</i>
$w_{\star} \in \mathbb{R}^{d}$ Teacher (true) weight vector $\hat{w} \in \mathbb{R}^{d}$ Learned weight vector (student) ε Adversarial perturbation budget $\ \cdot\ _{p}$ ℓ_{p} norm, defined as $\ x\ _{p} = (\sum_{i} x_{i} ^{p})^{1/p}$ p^{*} Dual exponent of p , satisfying $1/p + 1/p^{*} = 1$ $E_{rob}(\hat{w})$ Robust generalization error $E_{gen}(\hat{w})$ Standard generalization error $E_{bnd}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization strength parameter $\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_{w} Matrix defining Mahalanobis norm for regularization m^{*} Asymptotic overlap between teacher and student weights P^{*} Asymptotic dual norm of student weights P^{*} Asymptotic dual norm of student weights	$\star \in \mathbb{R}^d$ Teacher (true) weight vector $\in \mathbb{R}^d$ Learned weight vector (student)Adversarial perturbation budget ℓ_p norm, defined as $ x _p = (\sum_i x_i ^p)^{1/p}$ Dual exponent of p , satisfying $1/p + 1/p^* = 1$ rob (\hat{w}) Robust generalization errorgen (\hat{w}) Standard generalization errorbind (\hat{w}) Boundary error (difference between robust and standard error)Regularization strength parameter w)Regularization function \cdot)Surrogate loss function \cdots Output channel (conditional probability of labels)nInput distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights*Asymptotic dual norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table	$y_i \in \{-1, +1\}$	Binary label for sample <i>i</i>
$ \begin{split} \hat{w} \in \mathbb{R}^d & \text{Learned weight vector (student)} \\ \varepsilon & \text{Adversarial perturbation budget} \\ \ \cdot\ _p & \ell_p \text{ norm, defined as } \ x\ _p = (\sum_i x_i ^p)^{1/p} \\ p^* & \text{Dual exponent of } p, \text{ satisfying } 1/p + 1/p^* = 1 \\ E_{\text{rob}}(\hat{w}) & \text{Robust generalization error} \\ E_{\text{gen}}(\hat{w}) & \text{Standard generalization error} \\ E_{\text{bnd}}(\hat{w}) & \text{Boundary error (difference between robust and standard error)} \\ \lambda & \text{Regularization strength parameter} \\ \widetilde{r}(w) & \text{Regularization function} \\ g(\cdot) & \text{Surrogate loss function} \\ \mathbb{P}_{\text{out}} & \text{Output channel (conditional probability of labels)} \\ P_{\text{in}} & \text{Input distribution} \\ P_w & \text{Prior distribution on teacher weights} \\ \Sigma_{\delta} & \text{Matrix defining Mahalanobis norm for perturbations} \\ \Sigma_w & \text{Matrix defining Mahalanobis norm for regularization} \\ m^* & \text{Asymptotic overlap between teacher and student weights} \\ P^* & \text{Asymptotic dual norm of student weights} \\ \end{array} $	$\in \mathbb{R}^d$ Learned weight vector (student) Adversarial perturbation budget ℓ_p norm, defined as $ x _p = (\sum_i x_i ^p)^{1/p}$ Dual exponent of p , satisfying $1/p + 1/p^* = 1$ rob (\hat{w}) rob (\hat{w}) Robust generalization error Boundary error (difference between robust and standard error) Regularization strength parameter w) $w)$ Regularization function \cdot ω Prior distribution mut w w Prior distribution Matrix defining Mahalanobis norm for perturbations w ω Matrix defining Mahalanobis norm for regularization structure quare d ℓ_2 norm of student weights \star \star Asymptotic overlap between teacher and student weights \star \star Asymptotic dual norm of student weights \star Asymptotic dual norm of student weights \star Table 1: Notation Table	$oldsymbol{w}_{\star}\in\mathbb{R}^{d}$	Teacher (true) weight vector
$ \begin{split} \varepsilon & \text{Adversarial perturbation budget} \\ \ \cdot \ _{p} & \ell_{p} \text{ norm, defined as } \ x\ _{p} = (\sum_{i} x_{i} ^{p})^{1/p} \\ p^{*} & \text{Dual exponent of } p, \text{ satisfying } 1/p + 1/p^{*} = 1 \\ E_{\text{rob}}(\hat{w}) & \text{Robust generalization error} \\ E_{\text{gen}}(\hat{w}) & \text{Standard generalization error} \\ E_{\text{bnd}}(\hat{w}) & \text{Boundary error (difference between robust and standard error)} \\ \lambda & \text{Regularization strength parameter} \\ \widetilde{r}(w) & \text{Regularization function} \\ g(\cdot) & \text{Surrogate loss function} \\ \mathbb{P}_{\text{out}} & \text{Output channel (conditional probability of labels)} \\ P_{\text{in}} & \text{Input distribution} \\ P_{w} & \text{Prior distribution nt eacher weights} \\ \Sigma_{\delta} & \text{Matrix defining Mahalanobis norm for perturbations} \\ \Sigma_{w} & \text{Matrix defining Mahalanobis norm for regularization} \\ m^{*} & \text{Asymptotic overlap between teacher and student weights} \\ P^{*} & \text{Asymptotic dual norm of student weights} \\ P^{*} & \text{Asymptotic dual norm of student weights} \\ \end{array}$	Adversarial perturbation budget $\ _{p}$ ℓ_{p} norm, defined as $\ x\ _{p} = (\sum_{i} x_{i} ^{p})^{1/p}$ Dual exponent of p , satisfying $1/p + 1/p^{*} = 1$ $rob(\hat{w})$ Robust generalization error $gen(\hat{w})$ Standard generalization error $bnd(\hat{w})$ Boundary error (difference between robust and standard error) $Regularization strength parameter$ w)Regularization function \cdot)Surrogate loss function out Output channel (conditional probability of labels) n Input distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization \star Asymptotic overlap between teacher and student weights \star Asymptotic dual norm of student weights	$\hat{oldsymbol{w}} \in \mathbb{R}^d$	Learned weight vector (student)
$ \begin{array}{ll} \ \cdot \ _{p} & \ell_{p} \text{ norm, defined as } \ x \ _{p} = (\sum_{i} x_{i} ^{p})^{1/p} \\ p^{*} & \text{Dual exponent of } p, \text{ satisfying } 1/p + 1/p^{*} = 1 \\ E_{\text{rob}}(\hat{w}) & \text{Robust generalization error} \\ E_{\text{gen}}(\hat{w}) & \text{Standard generalization error} \\ E_{\text{bnd}}(\hat{w}) & \text{Boundary error (difference between robust and standard error)} \\ \lambda & \text{Regularization strength parameter} \\ \widetilde{r}(w) & \text{Regularization function} \\ g(\cdot) & \text{Surrogate loss function} \\ \mathbb{P}_{\text{out}} & \text{Output channel (conditional probability of labels)} \\ P_{\text{in}} & \text{Input distribution} \\ P_{w} & \text{Prior distribution on teacher weights} \\ \Sigma_{\delta} & \text{Matrix defining Mahalanobis norm for perturbations} \\ \Sigma_{w} & \text{Matrix defining Mahalanobis norm for regularization} \\ m^{*} & \text{Asymptotic overlap between teacher and student weights} \\ p^{*} & \text{Asymptotic dual norm of student weights} \\ Table 1: Notation Table \\ \end{array}$	$\begin{aligned} \ _{p} & \ell_{p} \text{ norm, defined as } \ x\ _{p} = (\sum_{i} x_{i} ^{p})^{1/p} \\ \text{Dual exponent of } p, \text{ satisfying } 1/p + 1/p^{*} = 1 \\ \text{rob}(\hat{w}) & \text{Robust generalization error} \\ \text{gen}(\hat{w}) & \text{Standard generalization error} \\ \text{Boundary error (difference between robust and standard error)} \\ \text{Regularization strength parameter} \\ w) & \text{Regularization function} \\ \cdot) & \text{Surrogate loss function} \\ \text{output channel (conditional probability of labels)} \\ \text{n} & \text{Input distribution} \\ w & \text{Prior distribution on teacher weights} \\ \delta & \text{Matrix defining Mahalanobis norm for perturbations} \\ w & \text{Matrix defining Mahalanobis norm for regularization} \\ \star & \text{Asymptotic overlap between teacher and student weights} \\ \star & \text{Asymptotic dual norm of student weights} \\ \text{Table 1: Notation Table} \\ \end{aligned}$	ε	Adversarial perturbation budget
p^* Dual exponent of p , satisfying $1/p + 1/p^* = 1$ $E_{rob}(\hat{w})$ Robust generalization error $E_{gen}(\hat{w})$ Standard generalization error $E_{bnd}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization strength parameter $\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights p^* Asymptotic dual norm of student weights P^* Asymptotic dual norm of student weightsTable 1: Notation Table	Dual exponent of p , satisfying $1/p + 1/p^* = 1$ rob(\hat{w})Robust generalization errorgen(\hat{w})Standard generalization errorbind(\hat{w})Boundary error (difference between robust and standard error)Regularization strength parameter w)Regularization function \cdot)Surrogate loss functionoutOutput channel (conditional probability of labels)nInput distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization \star Asymptotic overlap between teacher and student weights \star Asymptotic dual norm of student weights \star Asymptotic dual norm of student weights \star Table 1: Notation Table	$\ \cdot\ _p$	ℓ_p norm, defined as $ x _p = (\sum_i x_i ^p)^{1/p}$
$\begin{array}{lll} E_{\rm rob}(\hat{w}) & {\rm Robust\ generalization\ error} \\ E_{\rm gen}(\hat{w}) & {\rm Standard\ generalization\ error} \\ E_{\rm bnd}(\hat{w}) & {\rm Boundary\ error\ (difference\ between\ robust\ and\ standard\ error)} \\ \lambda & {\rm Regularization\ strength\ parameter} \\ \widetilde{r}(w) & {\rm Regularization\ strength\ parameter} \\ \widetilde{r}(w) & {\rm Regularization\ function} \\ g(\cdot) & {\rm Surrogate\ loss\ function} \\ \mathbb{P}_{\rm out} & {\rm Output\ channel\ (conditional\ probability\ of\ labels)} \\ P_{\rm in} & {\rm Input\ distribution} \\ P_{\rm w} & {\rm Prior\ distribution\ no\ teacher\ weights} \\ \Sigma_{\delta} & {\rm Matrix\ defining\ Mahalanobis\ norm\ for\ regularization} \\ \Sigma_w & {\rm Matrix\ defining\ Mahalanobis\ norm\ for\ regularization} \\ m^* & {\rm Asymptotic\ overlap\ between\ teacher\ and\ student\ weights} \\ P^* & {\rm Asymptotic\ dual\ norm\ of\ student\ weights} \\ Table\ 1:\ {\rm Notation\ Table} \end{array}$	rob(\hat{w})Robust generalization errorgen(\hat{w})Standard generalization errorbnd(\hat{w})Boundary error (difference between robust and standard error)Regularization strength parameterw)Regularization function \cdot)Surrogate loss functionoutOutput channel (conditional probability of labels)nInput distributionwPrior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbationswMatrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights*Asymptotic dual norm of student weights*Table 1: Notation Table	p^*	Dual exponent of p, satisfying $1/p + 1/p^* = 1$
$E_{\text{gen}}(\hat{w})$ Standard generalization error $E_{\text{bnd}}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization strength parameter $\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic dual norm of student weights P^* Table 1: Notation Table	(\hat{w}) Standard generalization error $bnd(\hat{w})$ Boundary error (difference between robust and standard error) $Regularization strength parameterw)Regularization function\cdot)Surrogate loss functionoutOutput channel (conditional probability of labels)nInput distributionwPrior distribution on teacher weights\deltaMatrix defining Mahalanobis norm for perturbationswMatrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights*Asymptotic dual norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table$	$E_{\rm rob}(\hat{\boldsymbol{w}})$	Robust generalization error
$E_{bnd}(\hat{w})$ Boundary error (difference between robust and standard error) λ Regularization strength parameter $\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic dual norm of student weights P^* Table 1: Notation Table	bindBoundary error (difference between robust and standard error) Regularization strength parameter w)Regularization function ψ)Regularization function ψ)Surrogate loss function ψ)Output channel (conditional probability of labels) n Input distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights*Asymptotic squared ℓ_2 norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table	$E_{\rm gen}(\hat{\boldsymbol{w}})$	Standard generalization error
λ Regularization strength parameter $\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic dual norm of student weights P^* Table 1: Notation Table	Regularization strength parameter w)Regularization function \cdot)Surrogate loss function out Output channel (conditional probability of labels) n Input distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights \star Asymptotic dual norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table	$E_{\rm bnd}(\hat{\boldsymbol{w}})$	Boundary error (difference between robust and standard error)
$\tilde{r}(w)$ Regularization function $g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weightsTable 1: Notation Table	w)Regularization function w)Surrogate loss function out Output channel (conditional probability of labels) n Input distribution w Prior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights δ Asymptotic squared ℓ_2 norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table	λ	Regularization strength parameter
$g(\cdot)$ Surrogate loss function \mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weightsTable 1: Notation Table	\cdot)Surrogate loss function \circ Output channel (conditional probability of labels) n n n w Prior distribution w δ Matrix defining Mahalanobis norm for perturbations w Matrix defining Mahalanobis norm for regularization \star Asymptotic overlap between teacher and student weights \star Asymptotic squared ℓ_2 norm of student weights \star Asymptotic dual norm of student weights \star Table 1: Notation Table	$\widetilde{r}(w)$	Regularization function
\mathbb{P}_{out} Output channel (conditional probability of labels) P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weightsTable 1: Notation Table	OutputOutput channel (conditional probability of labels)nInput distributionwPrior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbationswMatrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights*Asymptotic squared ℓ_2 norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table	$g(\cdot)$	Surrogate loss function
P_{in} Input distribution P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weightsTable 1: Notation Table	nInput distributionwPrior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbationswMatrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights*Asymptotic squared ℓ_2 norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table•the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c•throughout the paper and particularly in these proofs. The table includes both the	Pout	Output channel (conditional probability of labels)
P_w Prior distribution on teacher weights Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weightsTable 1: Notation Table	wPrior distribution on teacher weights δ Matrix defining Mahalanobis norm for perturbationswMatrix defining Mahalanobis norm for regularization*Asymptotic overlap between teacher and student weights*Asymptotic squared ℓ_2 norm of student weights*Asymptotic dual norm of student weights*Table 1: Notation Table	$P_{\rm in}$	Input distribution
Σ_{δ} Matrix defining Mahalanobis norm for perturbations Σ_{w} Matrix defining Mahalanobis norm for regularization m^{\star} Asymptotic overlap between teacher and student weights q^{\star} Asymptotic squared ℓ_2 norm of student weights P^{\star} Asymptotic dual norm of student weightsTable 1: Notation Table	 δ Matrix defining Mahalanobis norm for perturbations Matrix defining Mahalanobis norm for regularization * Asymptotic overlap between teacher and student weights Asymptotic squared l₂ norm of student weights * Asymptotic dual norm of student weights Table 1: Notation Table othe technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c d throughout the paper and particularly in these proofs. The table includes both th 	P_w	Prior distribution on teacher weights
Σ_w Matrix defining Mahalanobis norm for regularization m^* Asymptotic overlap between teacher and student weights q^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weights Table 1: Notation Table	w Matrix defining Mahalanobis norm for regularization * Asymptotic overlap between teacher and student weights * Asymptotic squared ℓ_2 norm of student weights * Asymptotic dual norm of student weights * Asymptotic dual norm of student weights Table 1: Notation Table o the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c d throughout the paper and particularly in these proofs. The table includes both the	Σ_{δ}	Matrix defining Mahalanobis norm for perturbations
m^* Asymptotic overlap between teacher and student weights q^* Asymptotic squared ℓ_2 norm of student weights P^* Asymptotic dual norm of student weightsTable 1: Notation Table	 * Asymptotic overlap between teacher and student weights Asymptotic squared l₂ norm of student weights * Asymptotic dual norm of student weights Table 1: Notation Table the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c d throughout the paper and particularly in these proofs. The table includes both the 	Σ_w	Matrix defining Mahalanobis norm for regularization
q^{\star} Asymptotic squared ℓ_2 norm of student weights P^{\star} Asymptotic dual norm of student weightsTable 1: Notation Table	 Asymptotic squared l₂ norm of student weights * Asymptotic dual norm of student weights Table 1: Notation Table the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c d throughout the paper and particularly in these proofs. The table includes both the 	m^{\star}	Asymptotic overlap between teacher and student weights
P* Asymptotic dual norm of student weights Table 1: Notation Table	* Asymptotic dual norm of student weights Table 1: Notation Table o the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c d throughout the paper and particularly in these proofs. The table includes both th	q^{\star}	Asymptotic squared ℓ_2 norm of student weights
Table 1: Notation Table	Table 1: Notation Table o the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c	P^{\star}	Asymptotic dual norm of student weights
Table 1: Notation Table	Table 1: Notation Table o the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c		-
	the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c		Table 1: Notation Table
	the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c		
· · · · · · · · · · · · · · · · · · ·	the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a c		
the the technical proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table T a c	d throughout the paper and particularly in these proofs. The table includes both the	nto the technical	proofs of Theorems 3.5, 3.7, 3.14 and 3.15, we provide in Table 1 a con
nd the more specialized symbols that appear in the asymptotic analysis. We have			changes symbols that appear in the asymptotic analysis. We have t

699

711 712

713

714

follow.

We now proceed with the proofs. The first theorem that will be crucial in our subsequent analysis is the Convex Gaussian MinMax Theorem (CGMT), a powerful tool in high-dimensional probability theory. The CGMT provides a connection between two seemingly unrelated optimization problems under Gaussian conditioning. Essentially, it allows us to study the properties of a complex primary optimization problem (PO) by examining a simpler auxiliary optimization problem (AO). This theorem is particularly valuable in our context as it enables us to transform intricate high-dimensional problems into more tractable lower-dimensional equivalents, significantly simplifying our analysis and leading to Theorems 3.7 and 3.15.

The CGMT states that under certain conditions, the probabilistic behavior of the primary optimization problem involving a
 Gaussian matrix is upper and lower bounded by the behavior of an auxiliary problem involving only Gaussian vectors. This
 powerful result allows us to derive tight probability bounds and asymptotic predictions for the high-dimensional estimation
 problems considered in this manuscript.

We state the theorem in full generality.

Theorem A.1 (CGMT (Gordon, 1988; Thrampoulidis et al., 2014)). Let $G \in \mathbb{R}^{m \times n}$ be an i.i.d. standard normal matrix and $g \in \mathbb{R}^m$, $h \in \mathbb{R}^n$ two i.i.d. standard normal vectors independent of one another. Let S_w , S_u be two compact sets such

that $S_{w} \subset \mathbb{R}^{n}$ and $S_{u} \subset \mathbb{R}^{n}$. Consider the two following optimization problems for any continuous ψ on $S_{w} \times S_{u}$

$$\mathbf{C}(\boldsymbol{G}) := \min_{\boldsymbol{w} \in \mathcal{S}_{\boldsymbol{w}}} \max_{\boldsymbol{u} \in \mathcal{S}_{\boldsymbol{u}}} \boldsymbol{u}^{\top} \boldsymbol{G} \boldsymbol{w} + \psi(\boldsymbol{w}, \boldsymbol{u})$$
(23)

$$\mathcal{C}(\mathbf{g}, \mathbf{h}) := \min_{\boldsymbol{w} \in \mathcal{S}_{\boldsymbol{w}}} \max_{\boldsymbol{u} \in \mathcal{S}_{\boldsymbol{u}}} \|\boldsymbol{w}\|_2 \mathbf{g}^\top \boldsymbol{u} + \|\boldsymbol{u}\|_2 \mathbf{h}^\top \boldsymbol{w} + \psi(\boldsymbol{w}, \boldsymbol{u})$$
(24)

Then the following hold

1. For all $c \in \mathbb{R}$ *we have*

$$\mathbb{P}(\mathbf{C}(\boldsymbol{G}) < c) \le 2\mathbb{P}(\mathcal{C}(\mathbf{g}, \mathbf{h}) \le c)$$
(25)

2. Further assume that S_w and S_u are convex sets, ψ is convex-concave on $S_w \times S_u$. Then for all $c \in \mathbb{R}$

$$\mathbb{P}(\mathbf{C}(\mathbf{G}) > c) \le 2\mathbb{P}(\mathcal{C}(\mathbf{g}, \mathbf{h}) \ge c)$$
(26)

In particular for all $\mu \in \mathbb{R}$, t > 0 we have $\mathbb{P}(|\mathbf{C}(\mathbf{G}) - \mu| > t) \le 2\mathbb{P}(|\mathcal{C}(\mathbf{g}, \mathbf{h}) - \mu| \ge t)$.

In our analysis, we will employ a version of the CGMT applied to a general class of generalized linear models, as proved by Loureiro et al. (2021).

A.1. Notations and Definitions

In this paper, we extensively employ the concepts of Moreau envelopes and proximal operators, pivotal elements in convex analysis frequently encountered in recent works on high-dimensional asymptotic of convex problems (Boyd & Vandenberghe, 2004; Parikh & Boyd, 2014). For an in-depth analysis of their properties, we refer the reader to the cited literature. Here, we briefly outline their definition and the main properties for context.

Definition A.2 (Moreau Envelope). Given a convex function $f : \mathbb{R}^n \to \mathbb{R}$ we define its Moreau envelope as being

$$\mathcal{M}_{Vf(\cdot)}(\boldsymbol{\omega}) = \min_{\boldsymbol{x}} \left[\frac{1}{2V} \|\boldsymbol{x} - \boldsymbol{\omega}\|_{2}^{2} + f(\boldsymbol{x}) \right].$$
(27)

where the Moreau envelope can be seen as a function $\mathcal{M}_{Vf(\cdot)} : \mathbb{R}^n \to \mathbb{R}$.

Definition A.3 (Proximal Operator). Given a convex function $f : \mathbb{R}^n \to \mathbb{R}$ we define its Proximal operator as being

$$\mathcal{P}_{Vf(\cdot)}(\boldsymbol{\omega}) = \operatorname*{arg\,min}_{\boldsymbol{x}} \left[\frac{1}{2V} \|\boldsymbol{x} - \boldsymbol{\omega}\|_{2}^{2} + f(\boldsymbol{x}) \right].$$
(28)

where the Proximal operator can be seen as a function $\mathcal{P}_{Vf(\cdot)} : \mathbb{R}^n \to \mathbb{R}^n$.

Theorem A.4 (Gradient of Moreau Envelope (Thrampoulidis et al., 2018), Lemma D1). Given a convex function $f : \mathbb{R}^n \to \mathbb{R}$, we denote its Moreau envelope by $\mathcal{M}_{Vf(\cdot)}(\cdot)$ and its Proximal operator as $\mathcal{P}_{Vf(\cdot)}(\cdot)$. Then, we have:

$$\nabla_{\boldsymbol{\omega}} \mathcal{M}_{Vf(\cdot)}(\boldsymbol{\omega}) = \frac{1}{V} \left(\boldsymbol{\omega} - \mathcal{P}_{Vf(\cdot)}(\boldsymbol{\omega}) \right) \,. \tag{29}$$

Additionally we will use the following two properties

$$\mathcal{M}_{Vf(\cdot+\boldsymbol{u})}(\boldsymbol{\omega}) = \mathcal{M}_{Vf(\cdot)}(\boldsymbol{\omega}+\boldsymbol{u}), \quad \mathcal{P}_{Vf(\cdot+\boldsymbol{u})}(\boldsymbol{\omega}) = \boldsymbol{u} + \mathcal{P}_{Vf(\cdot)}(\boldsymbol{\omega}+\boldsymbol{u}), \quad (30)$$

which are easy to show from a change of variables inside the minimization.

Definition A.5 (Dual of a Number). We define the dual of a number $a \ge 0$ as being a^* as the only number such that $1/a + 1/a^* = 1$.

A.2. Assumptions and Preliminary Discussion

We restate here all the assumptions that we make for the problem.

Assumption A.6 (Estimation from the dataset). Given a dataset \mathcal{D} made of n pairs of input outputs $\{(x_i, y_i)\}_{i=1}^n$, where $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$ we estimate the vector \hat{w} as being

$$\hat{\boldsymbol{w}} \in \underset{\boldsymbol{w} \in \mathbb{R}^d}{\arg\min} \sum_{i=1}^n \max_{\|\boldsymbol{\delta}_i\| \le \varepsilon} g\left(y_i \frac{\boldsymbol{w}^\top (\boldsymbol{x}_i + \boldsymbol{\delta}_i)}{\sqrt{d}} \right) + \lambda \widetilde{r}(\boldsymbol{w}) , \qquad (31)$$

where $g: \mathbb{R} \to \mathbb{R}$ is a convex non-increasing function, $\lambda \in [0, \infty)$ and $\tilde{r}: \mathbb{R}^d \to \mathbb{R}$ a convex regularization function.

Assumption A.7 (High-Dimensional Limit). We consider the proportional high-dimensional regime where both the number of training data and input dimension $n, d \to \infty$ at a fixed ratio $\alpha := n/d$.

Assumption A.8 (Regularization functions and Attack Norms considered). We consider consider two settings for the perturbation norm $\|\cdot\|$ and the regularization function r. For the first one, the regularization function and the attack norm ℓ_p norms, defined as

$$\|\boldsymbol{x}\|_{p} = \left(\sum_{i=1}^{n} |x_{i}|^{p}\right)^{1/p}$$
(32)

for $p \in (1, \infty]$. We will refer to the index of the regularization function as r and to the index of the norm inside the inner maximization as p and we define p^* as the dual number of p (definition A.5).

For the second case, both the regularization function and the attack norm are Mahalanobis norms, defined as

$$\|\boldsymbol{x}\|_{\boldsymbol{\Sigma}} = \sqrt{\boldsymbol{x}^{\top} \boldsymbol{\Sigma} \boldsymbol{x}} \tag{33}$$

for a positive definite matrix Σ . We refer to the index of the matrix of the regularization function as Σ_w and to the matrix of the norm inside the inner maximization as Σ_{δ} . In this case, we define p = r = 2 (in order to unify notations) and we will thus talk about $p^* = r^* = 2$.

This setting considers most of the losses used in machine learning setups for binary classification, *e.g.* logistic, hinge, exponential losses. We additionally remark that with the given choice of regularization the whole cost function is coercive.

Assumption A.9 (Scaling of Adversarial Norm Constraint). We suppose that the value of ε scales with the dimension d such that $\varepsilon \sqrt[p^*]{d} = O_d(1)$.

Assumption A.10 (Data Distribution). We consider two cases of data distribution. Both of them will rely on the following general generative process. For each $i \in [n]$, the covariates $x_i \in \mathbb{R}^d$ are drawn i.i.d. from a data distribution $P_{in}(x)$. Then, the corresponding y_i is sampled independently from the conditional distribution P_{out} . More succinctly, one can write the data distribution for a given pair (x, y) as

$$P(\boldsymbol{x}, \boldsymbol{y}) = \int_{\mathbb{R}^d} \mathrm{d}\boldsymbol{w}_{\star} \mathbb{P}_{\mathrm{out}}\left(\boldsymbol{y} \middle| \frac{\langle \boldsymbol{w}_{\star}, \boldsymbol{x} \rangle}{\sqrt{d}} \right) P_{\mathrm{in}}(\boldsymbol{x}) P_{\boldsymbol{w}}(\boldsymbol{w}_{\star}), \tag{34}$$

The target weight vector $w_{\star} \in \mathbb{R}^d$ is drawn from a prior probability distribution P_w .

¹⁵ Our two cases differentiate in the following way. For the first case, we consider $P_{in}(\boldsymbol{x}) = \mathcal{N}_{\boldsymbol{x}}(\boldsymbol{0}, \mathrm{Id}_d)$ and $P_{\boldsymbol{w}}$ which is ¹⁶ separable, i.e. $P_{\boldsymbol{w}}(\boldsymbol{w}) = \prod_{i=1}^{d} P_{w}(w_i)$ for a distribution P_{w} in \mathbb{R} with finite variance $\operatorname{Var}(P_w) = \rho < \infty$.

For the second case, we consider $P_{in}(x) = \mathcal{N}_x(0, \Sigma_x)$ and $P_w(w) = \mathcal{N}_w(0, \Sigma_\theta)$.

819 Assumption A.11 (Limiting Convergence of Spectral Values). We suppose that $\Sigma_x, \Sigma_\delta, \Sigma_\theta, \Sigma_w$ are simultaneously 820 diagonalisable. We call $\Sigma_x = S^T \operatorname{diag}(\omega_i)S$, $\Sigma_\delta = S^T \operatorname{diag}(\zeta_i)S$ and $\Sigma_w = S^T \operatorname{diag}(w_i)S$. We define $\overline{\theta} = S\Sigma_x^T w_* / \sqrt{\rho}$. 821 We assume that the empirical distributions of eigenvalues and the entries of $\overline{\theta}$ jointly converge to a probability distribution μ 822 as

$$\sum_{i=1}^{d} \delta(\bar{\theta}_i - \bar{\theta}) \delta(\omega_i - \omega) \delta(\zeta_i - \zeta) \delta(w_i - w) \to \mu.$$
(35)

A.3. Problem Simplification

Recall that we start from the following optimization problem:

$$\Phi_d = \min_{\boldsymbol{w} \in \mathbb{R}^d} \sum_{i=1}^n \max_{\|\boldsymbol{\delta}_i\| \le \varepsilon} g\left(y_i \frac{\boldsymbol{w}^\top(\boldsymbol{x}_i + \boldsymbol{\delta}_i)}{\sqrt{d}} \right) + \lambda r(\boldsymbol{w}),$$
(36)

where $r(\cdot)$ is a convex regularization function and $g(\cdot)$ is a non-increasing loss function. The non-increasing property of g allows us to simplify the inner maximization, leading to an equivalent formulation

$$\Phi_d = \min_{\boldsymbol{w} \in \mathbb{R}^d} \sum_{i=1}^n g\left(y_i \frac{\boldsymbol{w}^\top \boldsymbol{x}_i}{\sqrt{d}} - \frac{\varepsilon}{\sqrt{d}} \|\boldsymbol{w}\|_\star\right) + \lambda \widetilde{r}(\boldsymbol{w}).$$
(37)

To facilitate our analysis, we introduce auxiliary variables $P = \|\boldsymbol{w}\|_{\star}^{p^{\star}}/d$ and \hat{P} (the Lagrange parameter relative to this variable), which allow us to decouple the norm constraints. This leads to a min-max formulation

$$\Phi_{d} = \min_{\boldsymbol{w} \in \mathbb{R}^{d}, P} \max_{\hat{P}} \sum_{i=1}^{n} g\left(y_{i} \frac{\boldsymbol{w}^{\top} \boldsymbol{x}_{i}}{\sqrt{d}} - \frac{\varepsilon}{\sqrt[p^{\star}]{d}} \sqrt[p^{\star}]{\nabla} \right) + \lambda \widetilde{r}(\boldsymbol{w}) + \hat{P} \|\boldsymbol{w}\|_{\star}^{p^{\star}} - dP\hat{P},$$
(38)

where we switched the value of ε for its value without the scaling in d. This reformulation is what will allow us to apply the CGMT in subsequent steps.

It's worth noting the significance of the scaling for ε as detailed in Assumption A.9. In the high-dimensional limit $d \to \infty$, it's essential that all terms in Φ_d exhibit the same scaling with respect to d. This careful scaling ensures that our asymptotic analysis remains well-behaved and meaningful in the high-dimensional regime.

A.4. Scalarization and Application of CGMT

To facilitate our analysis, we further introduce effective regularization and loss functions, $\tilde{\tilde{r}}$ and \tilde{g} , respectively. These functions are defined as

$$\widetilde{g}(\boldsymbol{y},\boldsymbol{z}) = \sum_{i=1}^{n} g\left(y_i \boldsymbol{z}_i - \frac{\varepsilon}{\frac{p^*}{\sqrt{d}}} \sqrt[p^*]{P}\right), \quad \widetilde{\widetilde{r}}(\boldsymbol{w}) = \widetilde{r}(\boldsymbol{w}) + \hat{P} \|\boldsymbol{w}\|_{p^*}^{p^*}.$$
(39)

A crucial step in our analysis involves inverting the order of the min-max optimization. We can justify this operation by considering the minimization with respect to $w \in \mathbb{R}^d$ at fixed values of \hat{P} and P. This reordering is valid due to the convexity of our original problem. Specifically, the objective function is convex in w and concave in \hat{P} and P, and the constraint sets are convex. Under these conditions, we apply Sion's minimax theorem, which guarantees the existence of a saddle point and allows us to interchange the order of minimization and maximization without affecting the optimal value.

This reformulation enables us to directly apply (?)Lemma 11]loureiro2021learning. This lemma represents a meticulous application of Theorem A.1 to scenarios involving non-separable convex regularization and loss functions. The result is a lower-dimensional equivalent of our original high-dimensional minimization problem that represent the limiting behavior of the solution of the high-dimensional problem.

866 Consequently, our analysis now focuses on a low-dimensional functional, which takes the form 867

$$\widetilde{\Phi} = \min_{P,m,\eta,\tau_1} \max_{\hat{P},\kappa,\tau_2,\nu} \left[\frac{\kappa\tau_1}{2} - \alpha \mathcal{L}_g - \frac{\eta}{2\tau_2} \left(\nu^2 \rho + \kappa^2 \right) - \frac{\eta\tau_2}{2} - \mathcal{L}_{\widetilde{r}} + m\nu - P\hat{P} \right]$$
(40)

where we have restored the min max order of the problem.

In this expression, g and h are independent Gaussian vectors with i.i.d. standard normal components. The terms \mathcal{L}_g and $\mathcal{L}_{\tilde{r}}$ represent the scaled averages of Moreau Envelopes (eq. (27))

$$\mathcal{L}_{g} = \frac{1}{n} \mathbb{E} \left[\mathcal{M}_{\frac{\tau_{1}}{\kappa} \widetilde{g}(\boldsymbol{y},\cdot)} \left(\frac{m}{\sqrt{\rho}} \boldsymbol{s} + \eta \boldsymbol{h} \right) \right]$$
(41)

$$\mathcal{L}_{\widetilde{r}} = \frac{1}{d} \mathbb{E} \left[\mathcal{M}_{\frac{\eta}{\tau_2} \widetilde{\widetilde{r}}(\cdot)} \left(\frac{\eta}{\tau_2} (\kappa \boldsymbol{g} + \nu \boldsymbol{w}_\star) \right) \right]$$
(42)
879

The extremization problem in eq. (40) is related to the original optimization problem in eq. (36) as it can be thought as the leading part in the limit $n, d \rightarrow \infty$.

This dimensional reduction is the step that allows us to study the asymptotic properties of our original high-dimensional
 problem through a more tractable low-dimensional optimization and thus have in the end a low dimensional set of equations
 to study.

It's important to note that the optimization problem $\tilde{\Phi}$ is still implicitly defined in terms of the dimension d and, consequently, as a function of the sample size n. We introduce two variables

$$\boldsymbol{w}_{\mathrm{eq}} = \mathcal{P}_{\frac{\eta^*}{\tau_2^*} \widetilde{\tau}(.)} \left(\frac{\eta^*}{\tau_2^*} \left(\nu^* \mathbf{t} + \kappa^* \mathbf{g} \right) \right), \quad \boldsymbol{z}_{\mathrm{eq}} = \mathcal{P}_{\frac{\tau_1^*}{\kappa^*} \widetilde{g}(.,\mathbf{y})} \left(\frac{m^*}{\sqrt{\rho}} \mathbf{s} + \eta^* \mathbf{h} \right)$$
(43)

where $(\eta^{\star}, \tau_2^{\star}, P^{\star}, \hat{P}^{\star}, \kappa^{\star}, \nu^{\star}, m^{\star}, \tau_1^{\star})$ are the extremizer points of $\tilde{\Phi}$.

Building upon Loureiro et al. (2021, Theorem 5), we can establish a convergence result. Let \hat{w} be an optimal solution of the problem defined in eq. (36), and let $\hat{z} = \frac{1}{\sqrt{d}} X \hat{w}$. For any Lipschitz function $\varphi_1 : \mathbb{R}^d \to \mathbb{R}$, and any separable, pseudo-Lipschitz function $\varphi_2 : \mathbb{R}^n \to \mathbb{R}$, there exist constants $\epsilon, C, c > 0$ such that

$$\mathbb{P}\left(\left|\phi_{1}\left(\frac{\hat{\mathbf{w}}}{\sqrt{d}}\right) - \mathbb{E}\left[\phi_{1}\left(\frac{\boldsymbol{w}_{\mathrm{eq}}}{\sqrt{d}}\right)\right]\right| \geq \epsilon\right) \leq \frac{C}{\epsilon^{2}}e^{-cn\epsilon^{4}} \\
\mathbb{P}\left(\left|\phi_{2}\left(\frac{\hat{\mathbf{z}}}{\sqrt{n}}\right) - \mathbb{E}\left[\phi_{2}\left(\frac{\boldsymbol{z}_{\mathrm{eq}}}{\sqrt{n}}\right)\right]\right| \geq \epsilon\right) \leq \frac{C}{\epsilon^{2}}e^{-cn\epsilon^{4}}$$
(44)

It demonstrates that the limiting values of any function depending on \hat{w} and \hat{z} can be computed by taking the expectation of the same function evaluated at w_{eq} or z_{eq} , respectively. This convergence property allows us to translate results from our low-dimensional proxy problem back to the original high-dimensional setting with high probability.

912 A.5. Derivation of Saddle Point equations

We now want to show that extremizing the values of $m, \eta, \tau_1, P, \hat{P}, \nu, \tau_2, \kappa$ lead to the optimal value $\tilde{\Phi}$ of eq. (40). We are going to directly derive the saddle point equations and then argue that in the high-dimensional limit they become exactly the ones reported in the main text.

We obtain the first set of derivatives that depend only on the loss function and the channel part by taking the derivatives with respect to m, η, τ_1, P to obtain

$$\frac{\partial}{\partial m} : \nu = \alpha \frac{\kappa}{n\tau_1} \mathbb{E} \left[\left(\frac{m}{\eta \rho} \mathbf{h} - \frac{\mathbf{s}}{\sqrt{\rho}} \right)^\top \mathcal{P}_{\frac{\tau_1}{\kappa} \widetilde{g}(.,\mathbf{y})} \left(\frac{m}{\sqrt{\rho}} \mathbf{s} + \eta \mathbf{h} \right) \right]
\frac{\partial}{\partial \eta} : \tau_2 = \alpha \frac{\kappa}{\tau_1} \eta - \frac{\kappa \alpha}{\tau_1 n} \mathbb{E} \left[\mathbf{h}^\top \mathcal{P}_{\frac{\tau_1}{\kappa} \widetilde{g}(.,\mathbf{y})} \left(\frac{m}{\sqrt{\rho}} \mathbf{s} + \eta \mathbf{h} \right) \right]
\frac{\partial}{\partial \tau_1} : \frac{\tau_1^2}{2} = \frac{1}{2} \alpha \frac{1}{n} \mathbb{E} \left[\left\| \frac{m}{\sqrt{\rho}} \mathbf{s} + \eta \mathbf{h} - \mathcal{P}_{\frac{\tau_1}{\kappa} \widetilde{g}(.,y)} \left(\frac{m}{\sqrt{\rho}} \mathbf{s} + \eta \mathbf{h} \right) \right\|_2^2 \right]
\frac{\partial}{\partial P} : \hat{P} = \frac{\alpha}{n} \partial_P \mathbb{E} \left[\mathcal{M}_{\frac{\tau_1}{\kappa} \widetilde{g}(\mathbf{y}, \cdot)} \left(\frac{m}{\sqrt{\rho}} \mathbf{s} + \eta \mathbf{h} \right) \right]$$
(45)

By taking the derivatives with respect to the remaining variables $\kappa, \nu, \tau_2, \hat{P}$ we obtain a set of equations depending on

regularization and prior over the teacher weights

The rewriting of these equations in the desired form in Theorems 3.7 and 3.15 follows from the same considerations as in Loureiro et al. (2021, Appendix C.2).

To perform this rewriting the first ingredient we need is the following change of variables

$$m \leftarrow m, \qquad q \leftarrow \eta^2 + \frac{m^2}{\rho}, \qquad V \leftarrow \frac{\tau_1}{\kappa}, \qquad P \leftarrow P,$$

$$\hat{V} \leftarrow \frac{\tau_2}{n}, \qquad \hat{q} \leftarrow \kappa^2, \qquad \qquad \hat{m} \leftarrow \nu, \qquad \hat{P} \leftarrow \hat{P}.$$
(47)

and the use of Isserlis' theorem (Isserlis, 1918) to simplify the expectation where Gaussian g, h vectors are present.

A.5.1. REWRITING OF THE CHANNEL SADDLE POINTS

To obtain specifically the form implied in the main text we introduce

$$\mathcal{Z}_0(y,\omega,V) = \int \frac{\mathrm{d}x}{\sqrt{2\pi V}} e^{-\frac{1}{2V}(x-\omega)^2} \delta\left(y - f^0(x)\right) \,, \tag{48}$$

where this definition is equivalent to the one presented in eq. (10). The function \mathcal{Z}_0 can be interpreted as a partition function of the conditional distribution $\mathbb{P}_{\rm out}$ and contains all of the information about the label generating process.

A.5.2. Specialization of Prior Saddle Points for ℓ_p norms

In the case of ℓ_p norms, we can leverage the separable nature of the regularization to simplify our equations. The key insight here is that the proximal operator of a separable regularization is itself separable. This property allows us to treat each dimension independently, leading to a significant simplification of our high-dimensional problem.

First, due to the separability, all terms depending on the proximal of either \tilde{g} or \tilde{r} simplify the n or d at the denominator. This cancellation is crucial as it eliminates the explicit dependence on the problem dimension, allowing us to derive dimension-independent equations.

Next, we introduce

$$\mathcal{Z}_{w}(\gamma,\Lambda) = \int \mathrm{d}w \, P_{w}(w) e^{-\frac{\Lambda}{2}w^{2} + \gamma w},\tag{49}$$

which, in turn, leads in the form shown in eq. (9).

A.5.3. SPECIALIZATION OF PRIOR SADDLE POINTS FOR MAHALANOBIS NORMS

In the case of Mahalanobis norm, the form of the proximal of the effective regularization function is specifically

$$\mathcal{P}_{V\tilde{\widetilde{r}}(\cdot)}(\boldsymbol{\omega}) = \operatorname*{arg\,min}_{\boldsymbol{z}} \left[\lambda \boldsymbol{z}^{\top} \boldsymbol{\Sigma}_{\boldsymbol{w}} \boldsymbol{z} + \hat{P} \boldsymbol{z}^{\top} \boldsymbol{\Sigma}_{\boldsymbol{\delta}} \boldsymbol{z} + \frac{1}{2V} \|\boldsymbol{z} - \boldsymbol{\omega}\|_{2}^{2} \right] = \frac{1}{V} \left(2\hat{P} \boldsymbol{\Sigma}_{\boldsymbol{\delta}} + 2\lambda \boldsymbol{\Sigma}_{\boldsymbol{w}} + \frac{1}{V} \right)^{-1} \boldsymbol{\omega}$$
(50)

By substituting this explicit form into the equations from eq. (46), we obtain a set of simplified equations that still depends

990 on the dimension

991
992
$$m = \frac{1}{d} \operatorname{tr} \left[\hat{m} \boldsymbol{\Sigma}_{\boldsymbol{x}}^{\top} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^{\top} \boldsymbol{\Sigma}_{\boldsymbol{x}} \left(\lambda \boldsymbol{\Sigma}_{\boldsymbol{w}} + \hat{P} \boldsymbol{\Sigma}_{\boldsymbol{\delta}} + \hat{V} \boldsymbol{\Sigma}_{\boldsymbol{x}} \right)^{-1} \right]$$

 $q = \frac{1}{d} \operatorname{tr} \left[\left(\hat{m}^2 \boldsymbol{\Sigma}_{\boldsymbol{x}}^\top \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top \boldsymbol{\Sigma}_{\boldsymbol{x}} + \hat{q} \boldsymbol{\Sigma}_{\boldsymbol{x}} \right) \boldsymbol{\Sigma}_{\boldsymbol{x}} \left(\lambda \boldsymbol{\Sigma}_{\boldsymbol{w}} + \hat{P} \boldsymbol{\Sigma}_{\boldsymbol{\delta}} + \hat{V} \boldsymbol{\Sigma}_{\boldsymbol{x}} \right)^{-2} \right]$ $V = \frac{1}{d} \operatorname{tr} \left[\boldsymbol{\Sigma}_{\boldsymbol{x}} \Big(\lambda \boldsymbol{\Sigma}_{\boldsymbol{w}} + \hat{P} \boldsymbol{\Sigma}_{\boldsymbol{\delta}} + \hat{V} \boldsymbol{\Sigma}_{\boldsymbol{x}} \Big)^{-1} \right]$

996

1007

1012

1019

1022

1026

997 998 999

$$P = \frac{1}{d} \operatorname{tr} \left[\left(\hat{m}^2 \boldsymbol{\Sigma}_{\boldsymbol{x}}^\top \boldsymbol{\theta}_0 \boldsymbol{\theta}_0^\top \boldsymbol{\Sigma}_{\boldsymbol{x}} + \hat{q} \boldsymbol{\Sigma}_{\boldsymbol{x}} \right) \boldsymbol{\Sigma}_{\boldsymbol{\delta}} \left(\lambda \boldsymbol{\Sigma}_{\boldsymbol{w}} + \hat{P} \boldsymbol{\Sigma}_{\boldsymbol{\delta}} + \hat{V} \boldsymbol{\Sigma}_{\boldsymbol{x}} \right)^{-1} \right]$$

1000 The final step involves taking the high-dimensional limit of these equations. Here, we leverage our assumptions about the 1001 trace of the relevant matrices to further simplify the expressions so that they only depend on the limiting distribution μ from 1002 Assumption A.11. 1003

1004 Specifically, the assumptions on the trace allow us to replace certain high-dimensional operations with scalar quantities, 1005 effectively reducing the dimensionality of our problem. This dimensionality reduction is crucial for obtaining tractable 1006 equations in the high-dimensional limit. In the end we obtain the equations in eq. (16).

1008 A.6. Different channels and Prior functions 1009

We want to show how the different functions $\mathcal{Z}_0, \mathcal{Z}_w$ look like for some choices of output channel and prior in the data model. For the case of a probit output channel, we have by direct calculation

$$\mathcal{Z}_0(y,\omega,V) = \frac{1}{2}\operatorname{erfc}\left(-y\frac{\omega}{\sqrt{2(V+\tau^2)}}\right)$$
(52)

(51)

For the case of a channel of the form $y = \operatorname{sign}(z) + \sqrt{\Delta^* \xi}$, one has that 1016

$$\mathcal{Z}_{0}(y,\omega,V) = \mathcal{N}_{y}\left(1,\Delta^{\star}\right)\frac{1}{2}\left(1 + \operatorname{erf}\left(\frac{\omega}{\sqrt{2V}}\right)\right) + \mathcal{N}_{y}\left(-1,\Delta^{\star}\right)\frac{1}{2}\left(1 - \operatorname{erf}\left(\frac{\omega}{\sqrt{2V}}\right)\right)$$
(53)

For the choices of the prior over the teacher weights, we have for a Gaussian prior that

$$\mathcal{Z}_{\rm w}(\gamma,\Lambda) = \frac{1}{\sqrt{\Lambda+1}} e^{\gamma^2/2(\Lambda+1)}$$
(54)

or for sparse binary weights 1025

$$\mathcal{Z}_{\rm w}(\gamma,\Lambda) = \rho + e^{-\frac{\Lambda}{2}}(1-\rho)\cosh(\gamma)$$
(55)

A.7. Error Metrics 1028

To derive the form of the generalization error the procedure is the same as detailed in Aubin et al. (2020) or in Mignacco et al. (2020, Appendix A). We report here the final form being

$$E_{\rm gen} = \frac{1}{\pi} \arccos\left(\frac{m}{\sqrt{(\rho + \tau^2)q}}\right) \tag{56}$$

To derive the form for the boundary error one can proceed in the same way as Gerace et al. (2021, Appendix D) and obtain 1036

$$E_{\rm bnd} = \int_0^{\varepsilon \frac{p^* \sqrt{P}}{\sqrt{q}}} \operatorname{erfc}\left(-\frac{m}{\sqrt{q}}\lambda \frac{1}{\sqrt{2(\rho + \tau^2 - \frac{m^2}{q})}}\right) \frac{e^{-\frac{1}{2}\lambda^2}}{\sqrt{2\pi}} \,\mathrm{d}\lambda \tag{57}$$

We are also interested in the average teacher margin defined as

 $\mathbb{E}[y \boldsymbol{w}_{\star}^{\top} \boldsymbol{x}]$ (58)

On the Geometry of Regularization in Adversarial Training



Figure 5: Scaling of the overlap parameters in the low sample complexity regime for $p = \infty$, $\varepsilon = 0.3$, $\rho = 1$ and $\lambda = 10^{-3}$. The numbers presented in the legends are the linear fit in log-log scale of the dashed part.

¹⁰⁶³ which can be expressed as a function of the solutions of the saddle point equations as follows:

$$\sqrt{\frac{2}{\pi}} \frac{\sqrt{\rho}}{\sqrt{1 + \frac{\tau^2}{\rho}}} \tag{59}$$

1069 A.8. Asymptotic in the low sample complexity regime

This section examines the asymptotic behavior of our model in the regime of low sample complexity. Our analysis is motivated by numerical observations of the overlaps m, q, P, V in the small α regime, as illustrated in Figure 5.

¹⁰⁷³ Based on these observations, we propose a general scaling ansatz for the overlap parameters (solutions of the equations ¹⁰⁷⁴ presented in Theorems 3.7 and 3.15) as functions of the sample complexity α

$$m^{\star} = m_0 \alpha^{\delta_m} , \quad q^{\star} = q_0 \alpha^{\delta_q} , \quad V^{\star} = V_0 \alpha^{\delta_V} , \quad P^{\star} = P_0 \alpha^{\delta_P} , \tag{60}$$

1078 where the values with a zero subscript do not depend on α and the exponents are all positive. We focus on the noiseless case 1079 $\tau = 0$.

We are interested in the expansion of the generalization error and the boundary error, keeping only the most relevant terms in the limit $\alpha \to 0^+$. For the generalization error we have

$$E_{\rm gen} = \frac{1}{\pi} \arccos\left(\frac{m^{\star}}{\sqrt{\rho q^{\star}}}\right) = \frac{1}{2} - \frac{m_0}{\pi \sqrt{\rho q_0}} \alpha^{\delta_m - \frac{\delta_q}{2}} + o\left(\alpha^{\delta_m - \delta_q/2}\right) \tag{61}$$

⁶ and for the boundary error a similar expansion leads to

100

1088

1090

1083 1084 1085

1077

1062

$$E_{\rm bnd} = \int_0^{\frac{p_{\sqrt{P^\star}}}{q^\star}} \operatorname{erfc}\left(\frac{-\frac{m^\star}{\sqrt{q^\star}}\nu}{\sqrt{2(\rho - (m^\star)^2/q^\star)}}\right) \frac{e^{-\frac{\nu^2}{2}}}{\sqrt{2\pi}} \,\mathrm{d}\nu = \frac{\varepsilon_g \, {}^{p_{\sqrt{P_0}}}}{\sqrt{2\pi q_0}} \alpha^{\delta_P/p^\star - \delta_q/2} + \frac{\theta_0}{2\pi} \alpha^{2\delta_P/p^\star + \delta_m - 2\delta_q} + o(\alpha^\kappa) \quad (62)$$

1091 1092 where $\kappa = \max(\delta_P/p^* - \delta_q/2, 2\delta_P/p^* + \delta_m - 2\delta_q).$

¹⁰⁹³ Numerical simulations reveal a clear distinction in the low α regime between cases where the regularization parameter ¹⁰⁹⁴ $r = p^*$ and $r \neq p^*$. Figure 5 illustrates this difference for a fixed regularization parameter λ . We identify two scenarios that ¹⁰⁹⁵ characterize the behavior of the leading term in the boundary error expansion

1096 1097

When $\delta_P/p^* > \delta_q/2$: This occurs when $p^* = r = 1$. In this case, the leading term has a positive exponent, causing it to vanish as $\alpha \to 0$.

1100 When $\delta_P/p^* = \delta_q/2$: This scenario arises when $r \neq p^* = 1$. Here, the exponent of the leading term becomes zero, 1101 resulting in a constant term independent of α .

1103 Notably, in all cases we've examined, the second terms in both the generalization error and boundary error expansions
 1104 consistently approach zero in the limit of low sample complexity.
 1105

B. Rademacher Complexity Analysis

1108 Missing proofs for main text results.

Proposition B.1. Let $\varepsilon, \sigma > 0$. Consider a sample $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$, and let $\mathcal{H}_{\tilde{r}}$ be the hypothesis class defined in eq. (17). Then, it holds:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\widetilde{r}}) \leq \max_{i \in [n]} {}_{r} \|\mathbf{x}_{i}\|_{\star} \mathcal{W}_{\widetilde{r}} \sqrt{\frac{2}{\sigma n}} + \frac{\varepsilon}{2\sqrt{n}} \sup_{\boldsymbol{w}: \widetilde{r}(\boldsymbol{w}) \leq \mathcal{W}_{\widetilde{r}}^{2}} \|\boldsymbol{w}\|_{\star},$$
(63)

1115 where $_{r} \|\cdot\|_{\star}, \|\cdot\|_{\star}$ denote the dual norm of $_{r} \|\cdot\|, \|\cdot\|$, respectively.

1117 Proof. We have:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\widetilde{r}}) = \mathbb{E}_{\sigma} \left[\frac{1}{n} \sup_{h \in \mathcal{H}_{\widetilde{r}}} \sum_{i=1}^{n} \sigma_{i} \min_{\|\mathbf{x}_{i}' - \mathbf{x}_{i}\| \leq \varepsilon} y_{i} h(\mathbf{x}_{i}') \right] \\
= \mathbb{E}_{\sigma} \left[\frac{1}{n} \sup_{h \in \mathcal{H}_{\widetilde{r}}} \sum_{i=1}^{n} \sigma_{i} y_{i} \langle \boldsymbol{w}, \mathbf{x}_{i} \rangle - \varepsilon \| \boldsymbol{w} \|_{\star} \right] \qquad (\text{Def. of dual norm}) \\
\leq \mathbb{E}_{\sigma} \left[\frac{1}{n} \sup_{h \in \mathcal{H}_{\widetilde{r}}} \sum_{i=1}^{n} \sigma_{i} y_{i} \langle \boldsymbol{w}, \mathbf{x}_{i} \rangle \right] + \mathbb{E}_{\sigma} \left[\frac{1}{n} \sup_{h \in \mathcal{H}_{\widetilde{r}}} \sum_{i=1}^{n} -\varepsilon \sigma_{i} \| \boldsymbol{w} \|_{\star} \right] \qquad (\text{Subadditivity of supremum}) \\
= \hat{\mathfrak{R}}_{\mathrm{S}}(\mathcal{H}_{\widetilde{r}}) + \frac{\varepsilon}{2} \mathbb{E}_{\sigma} \left[\left| \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \right| \right] \sup_{\boldsymbol{w}: \widetilde{r}(\boldsymbol{w}) \leq \mathcal{W}_{\widetilde{r}}^{2}} \| \boldsymbol{w} \|_{\star}. \qquad (\text{Symmetry of } \sigma)$$

For the first term, i.e. the "clean" Rademacher Complexity, we plug in (Kakade et al., 2008, Theorem 1). By Jensen's inequality, we have for the second term:

$$\mathbb{E}_{\sigma}\left[\left|\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\right|\right] \leq \sqrt{\mathbb{E}_{\sigma}\left[\left(\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\right)^{2}\right]} = \frac{1}{\sqrt{n}},\tag{65}$$

1137 which concludes the proof.

Corollary B.2. Let $\varepsilon > 0$. Then:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\|\cdot\|_{2}^{2}}) \leq \frac{\max_{i \in [n]} \|\mathbf{x}_{i}\|_{2} \mathcal{W}_{2}}{\sqrt{m}} + \frac{\varepsilon \mathcal{W}_{2}}{2\sqrt{n}} \sqrt{\lambda_{\min}^{-1}(\boldsymbol{\Sigma}_{\boldsymbol{\delta}})}.$$
(66)

Proof. Leveraging Proposition 4.2, the first term of the RHS follows from the fact that the squared ℓ_2 norm is 1-strongly convex (w.r.t itself). For the second term, we have that the dual norm of $\|\cdot\|_{\Sigma_{\delta}}$ is given by $\|\cdot\|_{\Sigma_{\delta}^{-1}} = \sqrt{\langle w, \Sigma_{\delta}^{-1}w \rangle}$. Then, it holds:

$$\sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{2}^{2} \leq \mathcal{W}_{2}^{2}} \|\boldsymbol{w}\|_{\star} = \sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{2}^{2} \leq \mathcal{W}_{2}^{2}} \|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\delta}^{-1}}$$
$$= \mathcal{W}_{2} \sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{2} \leq 1} \|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\delta}^{-1}}$$
$$= \mathcal{W}_{2} \sqrt{\lambda_{\max}(\boldsymbol{\Sigma}_{\delta}^{-1})},$$
(67)

where the last equality follows from Courant–Fischer–Weyl's min-max principle.

Corollary B.3. Let $\Sigma_{w} = \sum_{i=1}^{d} \alpha_{i} \mathbf{v}_{i} \mathbf{v}_{i}^{T}$ and $\Sigma_{\delta} = \sum_{i=1}^{d} \lambda_{i} \mathbf{v}_{i} \mathbf{v}_{i}^{T}$, with $\mathbf{v}_{i} \in \mathbb{R}^{d}$ being orthonormal. Then:

 $\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{\|\cdot\|_{A}^{2}}) \leq \frac{\mathcal{W}_{A} \max_{i \in [n]} \|\mathbf{x}_{i}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}^{-1}}}{\sqrt{m}} + \frac{\varepsilon \mathcal{W}_{A}}{2\sqrt{n}} \sqrt{\max_{i \in [d]} \frac{1}{\lambda_{i} \alpha_{i}}}.$ (68)

Proof. For the worst-case part, we have:

$$\sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}}^{2} \leq \mathcal{W}_{A}^{2}} \|\boldsymbol{w}\|_{\boldsymbol{\Sigma}^{-1}} = \mathcal{W}_{A} \sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}} \leq 1} \sqrt{\langle \boldsymbol{w}, \boldsymbol{\Sigma}_{\delta}^{-1} \boldsymbol{w} \rangle}$$

$$= \mathcal{W}_{A} \sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}} \leq 1} \sqrt{\sum_{i=1}^{d} \lambda_{i}^{-1} \langle \boldsymbol{w}, \mathbf{v}_{i} \rangle^{2}}$$

$$= \mathcal{W}_{A} \sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}} \leq 1} \sqrt{\sum_{i=1}^{d} \frac{\lambda_{i}^{-1}}{\alpha_{i}} \alpha_{i} \langle \boldsymbol{w}, \mathbf{v}_{i} \rangle^{2}}$$

$$\leq \mathcal{W}_{A} \sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}} \leq 1} \sqrt{\max_{i\in[d]} \frac{\lambda_{i}^{-1}}{\alpha_{i}} \sum_{i=1}^{d} \alpha_{i} \langle \boldsymbol{w}, \mathbf{v}_{i} \rangle^{2}}$$

$$= \mathcal{W}_{A} \sup_{\boldsymbol{w}:\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}} \leq 1} \sqrt{\max_{i\in[d]} \frac{\lambda_{i}^{-1}}{\alpha_{i}}} \|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}} = \mathcal{W}_{A} \sqrt{\max_{i\in[d]} \frac{\lambda_{i}^{-1}}{\alpha_{i}}}.$$
(69)

1179 On the other hand, for $\boldsymbol{w} = \frac{1}{\sqrt{\alpha_j}} \mathbf{v}_j$ where $j \in \arg \max_{i \in [d]} \frac{\lambda_i^{-1}}{\alpha_i}$, it is $\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{w}}} = 1$ and also $\|\boldsymbol{w}\|_{\boldsymbol{\Sigma}_{\boldsymbol{\delta}}^{-1}} = \sqrt{\max_{i \in [d]} \frac{\lambda_i^{-1}}{\alpha_i}}$, 1180 so the above bound is tight.

B.1. Worst-case Rademacher complexity for ℓ_p norms

Awasthi et al. (2020) provides the following bound on the worst-case Rademacher complexity of linear hypothesis classes constrained in their ℓ_r norm.

Theorem B.4. Theorem 4 in (Awasthi et al., 2020) Let $\epsilon > 0$ and $p, r \ge 1$. Define $\mathcal{H}_{\tilde{r}} = \{\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle : \|\mathbf{w}\|_{r} \le 1\}$. Then, 1187 it holds:

$$\hat{\mathfrak{R}}_{\mathrm{S}}(\widetilde{\mathcal{H}}_{r}) \leq \hat{\mathfrak{R}}_{\mathrm{S}}(\mathcal{H}) + \epsilon \frac{\max(d^{1-\frac{1}{r}-\frac{1}{p}}, 1)}{2\sqrt{m}}.$$
(70)

1191 The bound above suggests regularizing the weights in the ℓ_r norm, $r = \frac{p}{p-1}$, for effectively controlling the estimation error 1192 of the class.

C. Parameter Exploration

This section presents the experimental details for all the figures in the main text and explore the model parameters in greater detail. For implementation details of our numerical procedures, please refer to Appendix D.

1199 C.1. Settings for Main Text Figures

All figures in the main text utilize the logistic loss function, defined as $g(x) = \log(1 + \exp(-x))$. Below, we detail the specific parameters for each figure.

Figure 1 We optimize the regularization parameter λ for each curve. Parameters: $\epsilon = 0.2$, noiseless regime ($\tau = 0$). Data points represent averages over 10 distinct data realizations with dimension d = 1000, varying sample size n to adjust α . Error bars indicate deviation from the mean.

Figure 2 Generated in the noiseless case ($\tau = 0$) with optimal regularization parameter λ . We optimize robust error for regularizations r = 2 and r = 1 independently, then compute their difference.



Figure 6: Robust error as a function of the regularization order r for two different p^* . By increasing the value of ε we have that the optimal value r^* gets close to p^* .



Figure 7: Robust error, generalization error and boundary error for different choices of regularization geometry r as a function of the sample complexity α . We see that the value of the errors increases with ε .

Figure 3 We employ a Strong Weak Feature Model (SWFM) as defined in Tanner et al. (2024). This model implements a block structure on all covariances (Σ_x , Σ_δ , Σ_θ , and Σ_w), with block sizes relative to dimension d denoted by ϕ_i for block i. We use two equal-sized feature blocks, totaling d = 1000. All matrices are block diagonal, with each block being diagonal. The values for each matrix are as follows

	Σ_x	Σ_δ	$\Sigma_{ heta}$	Case $\Sigma_w = \Sigma_\delta$	Case ℓ_2	Case $\Sigma_w = \Sigma_{\delta}^{-1}$
First Block	1	1	1	1	1	2.5
Second Block	1	2.5	1	2.5	1	1

All matrices are trace-normalized, with ε values as specified in the figure. Again error bars indicate the deviation from the mean.

Figure 4 We optimize the regularization parameter λ in the noiseless case ($\tau = 0$), with $\alpha = 1$. The inset is generated by conducting r sweeps for 10 distinct ε values. Each sweep comprises 50 points, with the minimum determined using np.argmin.

1258 C.2. Additional Parameter Exploration

1240

1251

1252

1256

We now present some additional exploration of the model in some different regimes.

Figure 6 These figures display theoretical results for attack perturbations constrained by ℓ_2 (Left) and $\ell_{3/2}$ (Right) norms. We vary ε as shown and use the noiseless regime ($\tau = 0$). Parameters: $\alpha = 0.1$, optimal λ . Each sweep comprises 50 points, with minima determined using np.argmin. Points on the curves indicate the minimum for each ε value.



1277 Figure 8: Robust generalization error as a function of regularization order r for fixed versus optimized regularization strength 1278 λ . The comparison illustrates that the impact of λ optimization does not change qualitatively the behavior of the optimal 1279 regularization geometry r^* as ε increases.

Figure 7 This figure illustrates generalization metrics as a function of α for various regularization geometries. We present results for two attack strengths: $\varepsilon = 0.1$ (Left) and $\varepsilon = 0.3$ (Right). Both use optimal λ values. This figure can be compared to Figure 1.

Figure 8 Both panels show robust generalization error versus regularization geometry r, with $\alpha = 0.1$. The right panel optimizes regularization strength λ , while the left uses a fixed value $\lambda = 10^{-4}$.

1288 1289 **D. Numerical Details**

The self-consistent equations from Theorems 3.7 and 3.15 are written in a way amenable to be solved via fixed-point iteration. Starting from a random initialization, we iterate through both the hat and non-hat variable equations until the maximum absolute difference between the order parameters in two successive iterations falls below a tolerance of 10^{-5} .

1294 To speed-up convergence we use a damping scheme, updating each order parameter at iteration *i*, designated as x_i , using 1295 $x_i := x_i \mu + x_{i-1}(1-\mu)$, with μ as the damping parameter.

¹²⁹⁶ ¹²⁹⁷Once convergence is achieved for fixed λ , hyper-parameters are optimized using a gradient-free numerical minimization ¹²⁹⁸procedure for a one dimensional minimization.

For each iteration, we evaluate the proximal operator numerically using SciPy's (Virtanen et al., 2020) Brent's algorithm for root finding (scipy.optimize.minimize_scalar). The numerical integration is handled with SciPy's quad method (scipy.integrate.quad), which provides adaptive quadrature of a given function over a specified interval. These numerical techniques allow us to evaluate the equations and perform the necessary integrations with the desired accuracy.

1303

1280 1281

1285

1286

1287

- 1304
- 1305
- 1306
- 1307 1308
- 1309
- 1310
- 311
- 1312
- 1314
- 1315
- 1316
- 1317
- 1318
- 1319