

PRECISE STATISTICAL ANALYSIS OF CLASSIFICATION ACCURACIES FOR ADVERSARIAL TRAINING

BY ADEL JAVANMARD^{*,†} AND MAHDI SOLTANOLKOTABI^{*,‡}

*University of Southern California**

Despite the wide empirical success of modern machine learning algorithms and models in a multitude of applications, they are known to be highly susceptible to seemingly small indiscernible perturbations to the input data known as *adversarial attacks*. A variety of recent adversarial training procedures have been proposed to remedy this issue. Despite the success of such procedures at increasing accuracy on adversarially perturbed inputs or *robust accuracy*, these techniques often reduce accuracy on natural unperturbed inputs or *standard accuracy*. Complicating matters further, the effect and trend of adversarial training procedures on standard and robust accuracy is rather counter intuitive and radically dependent on a variety of factors including the perceived form of the perturbation during training, size/quality of data, model overparameterization, etc. In this paper we focus on binary classification problems where the data is generated according to the mixture of two Gaussians with general anisotropic covariance matrices and derive a precise characterization of the standard and robust accuracy for a class of minimax adversarially trained models. We consider a general norm-based adversarial model, where the adversary can add perturbations of bounded ℓ_p norm to each input data, for an arbitrary $p \geq 1$. Our comprehensive analysis allows us to theoretically explain several intriguing empirical phenomena and provide a precise understanding of the role of different problem parameters on standard and robust accuracies.

1. Introduction Over the past decade there has been a tremendous increase in the use of machine learning models, and deep learning in particular, in a myriad of domains spanning computer vision and speech recognition, to robotics, healthcare and e-commerce. Despite wide empirical success in these and related domains, these modern learning models are known to be

[†]Supported in part by Sloan Research Fellowship in Mathematics, Google Faculty Research Award, Adobe Data Science Research Award and the NSF CAREER Award DMS-1844481.

[‡]Supported by the Packard Fellowship in Science and Engineering, Sloan Research Fellowship in Mathematics, an NSF-CAREER under award #1846369, the Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) under award #FA9550-18-1-0078, DARPA Learning with Less Labels (LwLL) and FastNICS programs, and NSF-CIF awards #1813877 and #2008443.

Keywords and phrases: precise high-dimensional asymptotics, adversarial training, binary classification

highly fragile and susceptible to *adversarial attacks*; even seemingly small imperceptible perturbations to the input data can significantly compromise their performance. As machine learning systems are increasingly being used in applications involving human subjects including healthcare and autonomous driving, such vulnerability can have catastrophic consequences. As a result there has been significant research over the past few years focused on proposing various *adversarial training* methods aimed at mitigating the effect of adversarial perturbations [20, 35, 42, 53, 68].

While adversarial training procedures have been successful in making machine learning models robust to adversarial attacks, their full effect on machine learning systems is not understood. Indeed, adversarial training procedures often behave in mysterious and somewhat counter intuitive ways. For instance, while they improve performance on adversarially perturbed inputs, this benefit often comes at the cost of decreasing accuracy on natural unperturbed inputs. This suggests that the two performance measures, *robust accuracy* –the accuracy on adversarially perturbed inputs– and the *standard accuracy* –accuracy on benign unperturbed inputs– may be fundamentally at conflict. Even more surprising, the performance of adversarial training procedure varies significantly in different settings. For instance, while adversarial trained models yield lower standard accuracy in comparison with non-adversarially trained counterparts, this behavior is completely reversed when there are very few training data with the standard accuracy of adversarially trained models outperforming that of non-adversarial models [65]. We refer the reader to Section 1.2 for a through discussion of recent empirical results that demonstrate how a variety of factors such as the adversary’s power, the size of training data, and model over-parameterization affect the performance of adversarially trained models.

To clearly demonstrate the surprising and counterintuitive behavior of adversarially trained models, we plot the behavior of such an approach in Figure 1. We consider a simple binary classification problem with the data generated according to a mixture of two isotropic Gaussians and depict the performance of a commonly used adversarial training procedure. In particular, in this figure, we plot the standard and robust accuracy of an adversarially trained linear classifier for different values of the adversary’s perceived power (measured in ℓ_∞ perturbations) and different sampling ratios (size of the training data divided by the number of parameters denoted by δ). We would like to highlight the highly non-trivial behavior of the standard and robust accuracy curves with respect to the adversary’s power and the sampling ratio. For instance, the standard accuracy first decreases, then increases and again decreases as a function of the adversary’s power. Furthermore, the

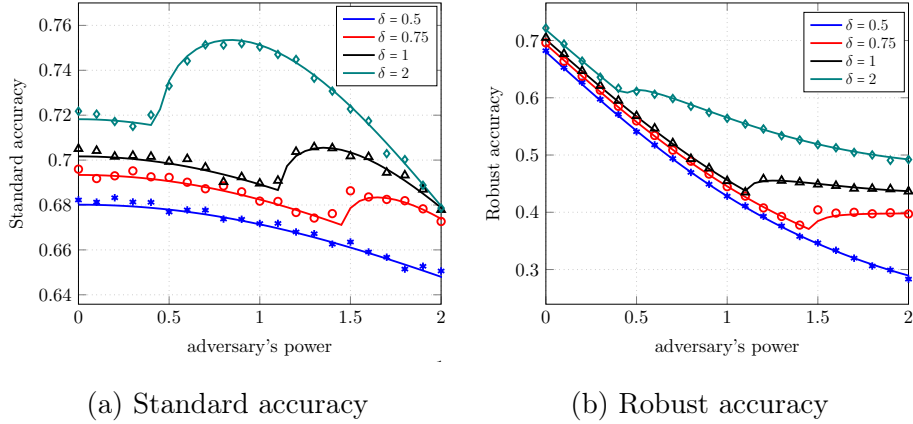


Fig 1: Depiction of standard and robust accuracies as a function of the adversary’s power with ℓ_∞ ($p = \infty$) perturbation for different values of δ (ratio of the size of the training data to the number of parameters in the model). Solid curves are theoretical predictions and dots are the empirical results. We refer to Figure 5 and Section 5.2 for further details.

exact nature of this curve is highly reliant on the sampling ratio δ . Similarly, for robust accuracy, we first observe a decreasing trend for all δ , but after some threshold depending on δ , robust accuracy increases and then decreases or stays constant. Even more surprising, as we will see in the forth-coming sections the behavior of these curves vary drastically for different forms of ℓ_p perturbations. This simple experiment clearly demonstrates the importance of having a precise theory for characterizing the rather nuanced performance of adversarial training procedures and demystify their behavior. Developing such a precise theoretical analysis is exactly the goal of this paper. Indeed, the solid curves in Figure 1 are based on our theoretical predictions!

1.1. Contributions In this paper we focus on binary classification problems where the data is generated according to the mixture of two Gaussians with general anisotropic covariance matrices and derive a precise characterization of the standard and robust accuracy for a class of minimax adversarially trained models. We consider a general norm-based adversarial model, where the adversary can add perturbations of bounded ℓ_p norm to each input data, for an arbitrary $p \geq 1$. We would like to emphasize that our theory provides a *precise characterization* of the performance of this class of adversarially trained models, rather than just upper bounds on the standard and robust

accuracies. Our analysis for such a broad setting allows us to capture several intriguing phenomena that we discuss next.

We show and theoretically prove an interesting phase transition phenomena holds for adversarial classification applied to the Gaussian mixture model. Specifically, we characterize a threshold δ_* for the ratio of size of training data to feature dimension, δ so that when $\delta < \delta_*$, the data is *robustly separable* with high probability, and for $\delta > \delta_*$ it is non-separable, with high probability. Here, *robust separability* is a generalization of the classical linear separability condition for data and roughly speaking means that there is a linear separator that correctly separates the two label classes with a positive margin that depends on the adversary's power. We precisely characterize the threshold δ_* in terms of various problem parameters including the mean and covariance of the mixture components, the adversary's power, and the ℓ_p perturbation norm. Interestingly, δ_* is related to the spherical width of a set defined in terms of the dual ℓ_q norm ($1/p + 1/q = 1$) conforming with classical notions of prior knowledge and complexity used in the compressive sensing literature.

Our precise theoretical characterization of standard and robust accuracies provides a precise understanding of the role that different problem parameters such as size/quality of the training data, feature covariates and means, model overparameterization ($1/\delta$), and the adversary's perceived power have during training on these performance measures. Surprisingly, our analysis reveals that the effects of these factors very much depend on the choice of perturbations norm ℓ_p . For example, in the robustly separable regime, we observe that for $p = 2$ adversarial training has no effect on standard accuracy, while for $p = 1$ and $p = \infty$ it hurts the standard accuracy. In the non-separable regime, we observe that for $p = 2$ adversarial training helps with improving the standard accuracy. However, for $p = \infty$ the adversarial training first improves the standard accuracy but as the training procedure hedges against stronger adversary, after some threshold on the adversary's power, we start to see a decrease in the standard accuracy of the resulting model. Interestingly, this threshold on the adversary's power varies with model overparameterization.

Lastly, a key ingredient of our analysis is a powerful extension of Gordon's Gaussian process inequality [21] known as the Convex Gaussian Minimax Theorem (CGMT) developed in [64] and further extended in [63, 12] for various learning settings. Using this technique we provide a precise prediction of the performance of adversarial training in terms of the optimal solutions to a convex-concave problem with a small number of scalar variables that can be easily solved by a low-dimensional gradient descent/ascent rather fast and accurately. In addition, this low-dimensional optimization problem can be significantly simplified for special cases of p (see Section 5 for details). While

CGMT has been used to study the behavior of regularized M-estimators, using this framework for the broad class of minimax adversarially trained models studied in this paper (including general anisotropic covariance matrices and general choice of ℓ_p norm for adversarial perturbations) poses significant technical challenges. Specifically, the intrinsic differences between ℓ_p geometries and the interaction between the class means the feature covariance matrix in the model requires a rather intricate and technical analysis.

1.2. *Related work* We briefly discuss the related literature along two lines.

Other models of adversarial perturbations. Another popular model for adversarial attacks on the models is the so-called distribution shifts, wherein the adversary can shift the test data distribution, making it different from the training distribution. The adversary is assumed to have limited manipulative power in terms of the Wasserstein distance between the test and the training distributions [58, 52, 44]. The articles [2, 44] study the robust loss $L(\boldsymbol{\theta}; \varepsilon) = \sup_{\nu \in B_\varepsilon(\mu)} \mathbb{E}_\nu[\ell(\mathbf{z}, \boldsymbol{\theta})]$, where $B_\varepsilon(\mu)$ is the ε ball around μ in the Wasserstein (W_p) distance for some $p \in [1, \infty)$, and the data $\mathbf{z} = (\mathbf{x}, y) \sim \mu$. A first order approximation of the robust loss $L(\boldsymbol{\theta}; \varepsilon)$ is given for small ε , in terms of a variation measure of the original loss ℓ . Such characterization is used in [44] to investigate the tradeoff between the standard and robust accuracies for various learning problems. Note that these work are focused on the population loss ($n \rightarrow \infty$, with d fixed). In comparison, in this paper we study norm bounded adversarial perturbations and work with empirical loss in asymptotic regime ($n, d \rightarrow \infty$, with $n/d = \delta$ fixed).

In adversarial training it is assumed that the modeler has access to clean (unperturbed) data and strives to construct a model that is resilient to potential adversarial perturbations of the test data. The article [36] considers a different adversarial setup in which an attacker can observe and modify all training data samples in an adversarial manner so as to maximize the estimation error caused by his attack. This work introduces the notion of adversarial influence function (AIF) to quantify the sensitivity of estimators to such adversarial attacks, and further derive the optimal estimator, among a certain class of estimator, that minimizes AIF.

Standard accuracy and robust accuracy tradeoffs. Several recent papers contain empirical results suggesting a potential trade-off between standard accuracy and robust accuracy. A few papers have started to shed light on the theoretical foundations of such tradeoffs [42, 57, 65, 54, 69, 30, 47, 14] often focusing on very specific models or settings. However, a comprehensive quantitative understanding of such tradeoffs is largely underdeveloped.

A central question we wish to address in this paper is whether there exists a fundamental conflict between robust accuracy and standard accuracy. We briefly mention a few papers that take a step towards addressing this question. In [65, 69], the authors provide examples of learning problems where no predictor can achieve both optimal standard accuracy and robust accuracy in the infinite data limit, pointing to such fundamental tradeoff. By contrast, [54] provides examples where there is no such tradeoff in the infinite data limit, in the sense that the optimal predictor performs well on both objectives, however a tradeoff is still observed with finite data. Despite this interesting progress a quantitative understanding of fundamental and algorithmic tradeoffs between standard and robust accuracies and how they are affected by various factors, such as overparameterization, adversary’s power and the data model is still missing. Such a result requires novel perspectives and analytical tools to precisely characterize the behavior of robust and standard accuracies, which is one of the motivating factors behind our current paper.

More closely related to this paper, in [30] the current authors used the convex Gaussian minimax framework to provide a precise characterization of standard and robust accuracies for linear regression, studying the fundamental conflict between these objectives along with algorithmic tradeoffs for specific minimax estimators. For classification problems, a recent paper [14] focuses on characterizing the optimal ℓ_2 and ℓ_∞ robust linear classifiers assuming access to the class means. This paper also studies some tradeoffs between standard and robust accuracies by contrasting this optimal robust classifier with the Bayes optimal classifier in a non-adversarial setting. This paper however does not directly study the tradeoffs of adversarial training procedures except for linear losses. A related publication [47] studies the generalization property of an adversarially trained model for classification on a Gaussian mixture model with a diagonal covariance matrix and a linear loss. In this setting, this work discusses the different effects that more training data can have on generalization based on the strength of the adversary. Using a linear loss in the above two classification papers is convenient as in this case the adversarially trained model admits a simple closed form representation. We also note that these two papers do not seem to focus on the high-dimensional regime where the number of training data grow in proportion to the number of parameters. In contrast, in this paper we focus on developing a comprehensive theory that provides a precise characterization of standard and robust accuracies and their tradeoffs in the high dimensional regime for a broad class of loss functions and covariance matrices. Such a comprehensive analysis allows us to better understand the role of the loss function in adversarial training. Indeed, as we demonstrate, the behavior of standard and robust accuracy for

nonlinear loss functions can be very different from linear losses. We also note that such a theoretical result requires much more intricate techniques as the adversarially trained model does not admit a simple closed form. Finally, we would like to note while in this paper we provide a precise understanding of the tradeoffs between standard and robust accuracies for commonly used adversarial training algorithms our work still does not address two tantalizing open questions: What is the optimal standard-robust accuracy tradeoff for a fixed ratio of sample size to dimension? Are there adversarial training approaches that achieve the optimal tradeoff between standard and robust accuracies universally over the range of adversary’s power.

2. Problem formulation In this section we discuss the problem setting and formulation of this paper in greater detail. After adopting some notations, we describe the adversarial training for binary classification in Section 2.1. Next, we discuss the data model and asymptotic setting studied in this paper in Section 2.2. Finally, in Section 2.3 we formally define the standard and robust classification accuracies in this model.

Notations. For a vector $\mathbf{v} \in \mathbb{R}^d$, we write $\|\mathbf{v}\|_{\ell_p}$ for the standard ℓ_p norm of \mathbf{v} , i.e., $\|\mathbf{v}\|_{\ell_p} = (\sum_i |v_i|^p)^{1/p}$. For a matrix $\mathbf{\Sigma}$, $\|\mathbf{\Sigma}\|$ indicates the spectrum norm of $\mathbf{\Sigma}$. Throughout, we say a probabilistic event holds ‘with high probability’, when its probability converges to one as $n \rightarrow \infty$. In addition, for a sequence of random variables $\{X_n\}_{n \in \mathbb{N}}$ and a constant c (independent of n) we write $\lim_{n \rightarrow \infty} X_n = c$, ‘in probability’ if $\forall \varepsilon > 0$ we have $\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - c| > \varepsilon) = 0$.

2.1. Adversarial training for binary classification In binary classification we have access to a training data set of n input-output pairs $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^d$ representing the input features and $y_i \in \{-1, +1\}$ representing the binary class label associated to each data point. Throughout we assume the data points (\mathbf{x}_i, y_i) are generated i.i.d. according to a distribution \mathbb{P} . To find a classifier that predicts the labels, one typically fits a function $f_{\boldsymbol{\theta}}$, parameterized by $\boldsymbol{\theta} \in \mathbb{R}^d$ to the training data via empirical risk minimization. In this paper we focus on linear classifiers of the form $f_{\boldsymbol{\theta}}(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\theta} \rangle$ in which case the training problem takes the form

$$(2.1) \quad \widehat{\boldsymbol{\theta}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i f_{\boldsymbol{\theta}}(\mathbf{x}_i)) = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle).$$

Here, ℓ is a loss and $\ell(y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle)$ approximately measuring the missclassification between the labels y_i and the output of the model $\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle$. Some common choices include logistic loss $\ell(t) = \log(1 + e^{-t})$, exponential loss $\ell(t) = e^{-t}$, and hinge loss $\ell(t) = \max(0, 1 - t)$. Once the parameter $\widehat{\boldsymbol{\theta}}$ is estimated one can

find the predicted label by simply calculating the sign of the model output $\widehat{y} = \text{sgn}(f_{\widehat{\theta}}(\mathbf{x})) = \text{sgn}(\langle \mathbf{x}, \widehat{\theta} \rangle)$.

Despite the widespread of empirical risk minimizers in supervised learning, these estimators are known to be highly vulnerable to even minute perturbations in the input features \mathbf{x}_i . In particular, it is known that even small, norm-bounded perturbations to the features that are imperceptible to the human eye, can lead to surprising miss-classification errors. These observations have spurred a surge of interest in adversarial training where the goal is to learn models that are robust against such adversarial perturbation. In this paper we focus on an adversarial training approach that is based on using a robust minimax loss [65, 42]. In our linear binary classification setting the robust minimax estimator takes the form

$$(2.2) \quad \widehat{\theta}^\varepsilon := \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max_{\|\delta_i\|_{\ell_p} \leq \varepsilon} \ell(y_i \langle \mathbf{x}_i + \delta_i, \theta \rangle).$$

The main intuition behind such an estimator is that although the learner has access to unperturbed training data, instead of fitting to that data she imitates potential adversarial perturbations to test data in the training data and aims to learn a model that performs well in the presence of such perturbations. One can also view this adversarial training approach as an implicit smoothing that tries to fit the same label y_i to all the features in the ε -neighborhood of \mathbf{x}_i simultaneously.

In this paper we focus on convex and decreasing losses such as the aforementioned logistic, exponential, and hinge losses. In such cases the inner maximization in (2.2) can be solved in closed form. In particular, the worst perturbation δ_i in terms of loss value is given by $\delta_i^* = \arg \min \{y_i \langle \delta_i, \theta \rangle : \|\delta_i\|_{\ell_p} \leq \varepsilon\}$, which by using Holder's inequality results in $y_i \langle \delta_i^*, \theta \rangle = -\varepsilon \|\theta\|_{\ell_q}$. Therefore the adversarially trained model $\widehat{\theta}^\varepsilon$ can be equivalently written as

$$(2.3) \quad \widehat{\theta}^\varepsilon := \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i \langle \mathbf{x}_i, \theta \rangle - \varepsilon \|\theta\|_{\ell_q}).$$

2.2. Data model and asymptotic setting We consider supervised binary classification under a Gaussian Mixture data Model (GMM). Concretely, each data point belongs to one of two classes $\{\pm 1\}$ with corresponding probabilities π_+ , π_- , so that $\pi_+ + \pi_- = 1$. Given the label $y_i \in \{-1, +1\}$ for data point i , the associated input/feature vectors $\mathbf{x}_i \in \mathbb{R}^d$ are generated independently according to the distribution $\mathbf{x}_i \sim \mathcal{N}(y_i \boldsymbol{\mu}, \boldsymbol{\Sigma})$, conditioned on y_i , where $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$. In other words the mean of feature vectors are $\pm \boldsymbol{\mu}$ depending on its class, and $\boldsymbol{\Sigma}$ is the covariance of features. We depict this mixture

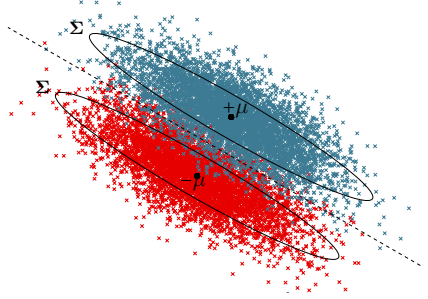


Fig 2: Depiction of the Mixture of Gaussian data model.

model in Figure 2. We next describe the asymptotic regime of interest and our assumptions in this paper.

Assumption 1 (Asymptotic Setting) *We focus on the following asymptotic regime:*

- (a) (Scaling of dimensions) $n \rightarrow \infty$ and $\frac{n}{d} \rightarrow \delta \in (0, \infty)$.
- (b) (Scaling of signal to noise ratio) We have $C_{\min} \leq \frac{\|\mu\|_{\ell_2}}{\|\Sigma\|} \leq C_{\max}$ for some positive constants C_{\min} and C_{\max} , which are independent of n and d .
- (c) (Scaling of adversary's power) We have $\varepsilon = \varepsilon_0 \|\mu\|_{\ell_p}$ for a constant $\varepsilon_0 \geq 0$ which we refer to as adversary's normalized power.

Assumption 1 (a) details our high-dimensional regime where the size of the training data n and the dimension of the features d grow proportionally with their ratio fixed at δ . We would like to note that while we focus on this asymptotic regime our theoretical technique can also demonstrate very accurate concentration around this asymptotic behavior. Assumption 1 (b) demonstrates the scaling of the signal to noise ratio and ensures that the distance between the centers of the two components $2\|\mu\|_{\ell_2}$ ('signal') is comparable to the projection of noise in any direction (noise). Finally, Assumption 1 details our scaling of the adversary's power. This scaling is justified as if the adversary could perturb data points \mathbf{x}_i by 2μ , she can flip the label of every data point, so that the learner cannot do better than random guessing. Since the perturbations can be chosen arbitrary from an ℓ_p ball of radius ε , we require ε to be comparable to $\|\mu\|_{\ell_p}$.

2.3. Standard and robust accuracies Our goal in this paper is to precisely characterize performance of the estimator $\hat{\theta}^\varepsilon$ in terms of two accuracies and understand the interplay between them. The two accuracies are *standard accuracy* which is the accuracy on unperturbed test data, and *robust accuracy*

which is the accuracy on adversarially perturbed test data. More formally *standard accuracy* quantifies the accuracy of an estimator on an unperturbed test data that is generated from the same distribution as the training data:

$$(2.4) \quad \text{SA}(\widehat{\boldsymbol{\theta}}) := \mathbb{P}\{\widehat{y} = y\}, \quad \text{where } (\mathbf{x}, y) \sim \mathbb{P}$$

Our second accuracy, called *robust accuracy* quantifies robustness of an estimator to adversarial perturbations in the test data. Specifically,

$$(2.5) \quad \text{RA}(\widehat{\boldsymbol{\theta}}) := \mathbb{E} \left[\min_{\|\boldsymbol{\delta}\|_{\ell_p} \leq \varepsilon} \mathbb{1}_{\{y\langle \mathbf{x} + \boldsymbol{\delta}, \widehat{\boldsymbol{\theta}} \rangle \geq 0\}} \right], \quad \text{where } (\mathbf{x}, y) \sim \mathbb{P}.$$

We end this section by stating a lemma that characterizes $\text{SA}(\widehat{\boldsymbol{\theta}})$ and $\text{RA}(\widehat{\boldsymbol{\theta}})$ under the Gaussian mixture model. We defer the proof to Appendix E.1.

Lemma 2.1 *Consider mixtures of Gaussian data model where $y_i \in \{-1, +1\}$ with corresponding probabilities π_-, π_+ and the feature vector distributed as $\mathbf{x} \sim \mathcal{N}(y\boldsymbol{\mu}, \boldsymbol{\Sigma})$, conditioned on y , where $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$. Then,*

$$(2.6) \quad \text{SA}(\widehat{\boldsymbol{\theta}}) := \Phi \left(\frac{\langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}} \rangle}{\|\boldsymbol{\Sigma}^{1/2} \widehat{\boldsymbol{\theta}}\|_{\ell_2}} \right),$$

$$(2.7) \quad \text{RA}(\widehat{\boldsymbol{\theta}}) := \Phi \left(\frac{\langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}} \rangle - \varepsilon \|\widehat{\boldsymbol{\theta}}\|_{\ell_q}}{\|\boldsymbol{\Sigma}^{1/2} \widehat{\boldsymbol{\theta}}\|_{\ell_2}} \right).$$

Here, $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$ is the cdf of a standard Gaussian distribution and q is such that $\frac{1}{p} + \frac{1}{q} = 1$.

By Lemma 2.1, characterizing $\text{SA}(\widehat{\boldsymbol{\theta}}^\varepsilon)$ and $\text{RA}(\widehat{\boldsymbol{\theta}}^\varepsilon)$ amounts to characterizing $\langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}}^\varepsilon \rangle$, $\|\boldsymbol{\Sigma}^{1/2} \widehat{\boldsymbol{\theta}}\|_{\ell_2}$, $\|\widehat{\boldsymbol{\theta}}^\varepsilon\|_{\ell_q}$, which constitutes the bulk of our analysis.

3. Prelude: two regimes for adversarial training Similar to normal classification, an interesting phenomena that arises in adversarial classification is that depending on the size of the training data there are two different regimes of operation: Robustly separable and non-separable. In the *robustly separable* regime there is a robust classifier that perfectly separates the training data, with a positive margin that depends on the adversary's power, while this is not possible in the non-separable case. We formally define this notion of robust separability below.

Definition 3.1 (Robust linear separability) *Given $\varepsilon > 0$ and $q \geq 1$, we call a training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, (ε, q) -separable if*

$$(3.1) \quad \exists \boldsymbol{\theta} \in \mathbb{R}^d : \quad \forall i \in [n], \quad y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} > 0.$$

We note that our notion of robust separability is closely related to the standard notion of separability by a linear classifier. In particular, using a simple rescaling argument¹ one can rewrite condition 3.1 as follows

$$(3.2) \quad \exists \boldsymbol{\theta}, \|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon} : \quad \forall i \in [n], \quad y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle > 1.$$

Therefore, robust separability is akin to linear separability of the data but with a budget constraint on the ℓ_q norm of the coefficients of the classifier.

When the training data is (ε, q) -separable (with ℓ_q the dual norm of ℓ_p), then the minimax estimator $\widehat{\boldsymbol{\theta}}^\varepsilon$ becomes unbounded and achieves zero adversarial training loss in (2.3). In other words, one can completely interpolate the data. This is due to the fact that if $\boldsymbol{\theta}$ is an (ε, q) -separator, then $c\boldsymbol{\theta}$ with $c \rightarrow \infty$ leads to zero adversarial training loss and since the loss is nonnegative it is optimal. Although the norm of $\widehat{\boldsymbol{\theta}}^\varepsilon$ tends to infinity in the separable regime, what matters for our linear classifier is the direction of $\widehat{\boldsymbol{\theta}}^\varepsilon$. However, in this separable regime even the direction of the optimal solution $(\frac{\widehat{\boldsymbol{\theta}}^\varepsilon}{\|\widehat{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2}})$ may not be unique. Even though there may be multiple optimal directions it is possible to show that the direction that gradient descent converges to is a specific maximum margin classifier. We formally state this result which is essentially a direct consequence of [41, 31] below.

Proposition 3.2 *Consider the adversarial training loss*

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell \left(y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right),$$

with the loss $\ell(t)$ obeying certain technical assumptions² which are satisfied for common classification losses such as logistic, exponential, and hinge losses. Then, the gradient descent iterates

$$\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_\tau - \mu \nabla \mathcal{L}(\boldsymbol{\theta}_\tau)$$

with a sufficiently small step size μ obey

$$(3.3) \quad \lim_{t \rightarrow \infty} \left\| \frac{\boldsymbol{\theta}_t}{\|\boldsymbol{\theta}_t\|_{\ell_2}} - \frac{\widetilde{\boldsymbol{\theta}}^\varepsilon}{\|\widetilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2}} \right\|_{\ell_2} = 0,$$

¹(3.1) \Rightarrow (3.2): Scaling by $\frac{1}{\varepsilon \|\boldsymbol{\theta}\|_{\ell_q}}$ we see that $y_i \langle \mathbf{x}_i, \tilde{\boldsymbol{\theta}} \rangle > 1$ for $\tilde{\boldsymbol{\theta}} = \frac{\boldsymbol{\theta}}{\varepsilon \|\boldsymbol{\theta}\|_{\ell_q}}$, and $\|\tilde{\boldsymbol{\theta}}\|_{\ell_q} = \frac{1}{\varepsilon}$ by definition. (3.2) \Rightarrow (3.1): Letting $c = \frac{1}{\varepsilon \|\boldsymbol{\theta}\|_{\ell_q}} \geq 1$, we also have $y_i \langle \mathbf{x}_i, c\boldsymbol{\theta} \rangle > 1$. Substituting for c and rearranging the terms we get (3.1).

²See [41, Assumption S3 in Appendix F]. We list these assumptions in Appendix F for readers' convenience.

where $\tilde{\boldsymbol{\theta}}^\varepsilon$ is the solution to the following max-margin problem

$$(3.4) \quad \begin{aligned} \tilde{\boldsymbol{\theta}}^\varepsilon = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \quad & \|\boldsymbol{\theta}\|_{\ell_2}^2 \\ \text{subject to} \quad & y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \geq 1. \end{aligned}$$

In the non-separable regime, as we show in the proof of Theorem 4.5 the minimizer $\tilde{\boldsymbol{\theta}}^\varepsilon$ is bounded. Moreover, the loss (2.3) is convex as it is pointwise maximum of a set of convex functions (see (2.2) and recall convexity of loss ℓ). Therefore, a variety of iterative methods (including gradient descent) can be used to converge to a global minimizer of (2.2). Theorem 4.5 also shows that all global minimizers of (2.2) have the same standard and robust accuracy.

4. Main results for isotropic features In this section we present our main results. For the sake of exposition, in this section we state our results for the case where the features are isotropic (i.e. $\boldsymbol{\Sigma} = \mathbf{I}$). We discuss our more general results with anisotropic features in Section 6. In this paper, we establish a sharp phase-transition characterizing the separability of the training data generated according to a Gaussian mixture model. Specifically, in our asymptotic regime (see Section 2.1) we characterize a threshold δ_* such that for $\delta < \delta_*$ the data is (ε, q) -separable, with high probability, and for $\delta > \delta_*$ it is non-separable, with high probability. This phase transition for robust separability is discussed in Section 4.1. We also precisely characterize the standard accuracy $\text{SA}(\tilde{\boldsymbol{\theta}}^\varepsilon)$ and the robust accuracy $\text{RA}(\tilde{\boldsymbol{\theta}}^\varepsilon)$ of the point that gradient descent converges to in both the separable and non-separable data regimes which are the subject of Sections 4.2 and 4.3, respectively. We then discuss the implications of our main results for the special cases of ℓ_p perturbations with $p = 1$, $p = 2$, and $p = \infty$ in Section 5.

4.1. Phase transition for robust data separability In this section we discuss our results for characterizing the phase transition for (ε, q) -separability under the Gaussian mixtures model. As detailed earlier in Section 2.2, in our asymptotic setting the dimension of the mean vector $\boldsymbol{\mu}$ (d) as well as the size of the training data (n) grow to infinity in proportion with each other $n/d = \delta$. To state our main result we need a few technical assumptions on the limiting behavior of the mean vector. We begin with a simple assumption on the convergence of the Euclidean norm of the mean vector.

Assumption 2 (Convergence of Euclidean norm of $\boldsymbol{\mu}$) *We assume the Euclidean norm of the mean vector converges to a bounded quantity, that is $\|\boldsymbol{\mu}\|_{\ell_2} \rightarrow V < \infty$, as $n \rightarrow \infty$ and $n/d \rightarrow \delta$.*

We note that for the isotropic case, the boundedness condition in Assumption 2 is already implied by Assumption 1(b).

Naturally, the separability threshold depends on the mean vector and the adversary's power. For instance, intuitively, one expects the separability threshold to decrease as the adversary's power or the length of the mean vector increases. We also expect the direction of the mean vector $\frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_{\ell_2}}$ to play a role. We capture these effects via the spherical width of a suitable set. Recall that the *spherical width* of a set $\mathcal{S} \subset \mathbb{R}^d$ is a measure of its complexity and is defined as

$$\omega_s(\mathcal{S}) = \mathbb{E} \left[\sup_{\mathbf{z} \in \mathcal{S}} \mathbf{z}^T \mathbf{u} \right],$$

where $\mathbf{u} \in \mathcal{S}^{d-1}$ is a vector chosen uniformly at random from the unit sphere. In particular, the appropriate set for characterizing the separability threshold takes the form

$$(4.1) \quad \mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu}) := \left\{ \mathbf{z} \in \mathbb{R}^d : \mathbf{z}^T \boldsymbol{\mu} = 0, \|\mathbf{z}\|_{\ell_2} \leq \alpha, \left\| \mathbf{z} + \theta \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_{\ell_2}} \right\|_{\ell_q} \leq \frac{1}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}} \right\},$$

where ε_0 is the adversary's scaled power per Assumption 1(c). Next assumption focuses on the spherical width convergence in our asymptotic regime.

Assumption 3 (Convergence of spherical width) *We assume the following limit exists*

$$(4.2) \quad \omega(\alpha, \theta, \varepsilon_0) := \lim_{n \rightarrow \infty} \omega_s(\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})).$$

As it will become clear later on in this section Assumptions 2 and 3 are trivially satisfied in various settings. With these assumptions in place we are ready to state our result precisely characterizing the separability threshold.

Theorem 4.1 *Consider a data set generated i.i.d. according to an isotropic Gaussian mixture data model per Section 2.2 and suppose the mean vector $\boldsymbol{\mu}$ obeys Assumptions 2 and 3. Also define*

$$(4.3) \quad \delta_* := \min_{\alpha \geq 0, \theta} \frac{\omega(\alpha, \theta, \varepsilon_0)^2}{\mathbb{E} \left[\left(1 - V\theta + \sqrt{\alpha^2 + \theta^2} g \right)_+^2 \right]},$$

where the expectation is taken with respect to $g \sim \mathcal{N}(0, 1)$. Then, under the asymptotic setting of Assumption 1, for $\delta < \delta_*$ the data are (ε, q) -separable

with high probability and for $\delta > \delta_*$, the data are non-separable, with high probability. Namely,

$$\begin{aligned}\delta < \delta_* &\Rightarrow \lim_{n \rightarrow \infty} \mathbb{P}(\text{data is } (\varepsilon, q)\text{-separable}) = 1, \\ \delta > \delta_* &\Rightarrow \lim_{n \rightarrow \infty} \mathbb{P}(\text{data is } (\varepsilon, q)\text{-separable}) = 0.\end{aligned}$$

Theorem 4.1 above precisely characterizes the separability threshold as a function of the adversary's power as well as properties of the mean vector. In particular since ω decreases with the increase in ε_0 , this theorem indicates that the separability threshold decreases as the adversary's power increases. This of course conforms with our natural intuition and is consistent with characterization (3.2). To better understand the implications of Theorem 4.1 we now consider some special cases.

- **Example 1 (Non-adversarial setting).** Our first example focuses on the non-adversarial setting where $\varepsilon_0 = 0$. In this case the ℓ_q constraint in definition of \mathcal{S} , given by (4.1), is void and the set \mathcal{S} becomes the intersection of ℓ_2 ball of radius α with the hyperplane of dimension $d - 1$ that is orthogonal to $\boldsymbol{\mu}$. Therefore $\omega(\alpha, \theta, \varepsilon_0) = \omega_s(\mathcal{S}) = \alpha$ and the separability threshold reduces to

$$\delta_* := \max_{\alpha \geq 0, \theta} \frac{\alpha^2}{\mathbb{E} \left[\left(1 - V\theta + \sqrt{\alpha^2 + \theta^2} g \right)_+^2 \right]}.$$

By the change of variables $(\alpha, \frac{\theta}{\alpha}) \rightarrow (\alpha, \theta)$, it is straightforward to see that optimal α is at $+\infty$ and the separability condition reduces to

$$\delta_* := \left(\min_{\theta} \mathbb{E} \left[\left(-V\theta + \sqrt{1 + \theta^2} g \right)_+^2 \right] \right)^{-1}.$$

- **Example 2 (ℓ_2 perturbation).** When $p = q = 2$, the set \mathcal{S} becomes the intersection of ℓ_2 ball of radius

$$R := \min \left(\alpha, \sqrt{\frac{1}{\varepsilon_0^2 \|\boldsymbol{\mu}\|_{\ell_2}^2} - \theta^2} \right),$$

with the hyperplane of dimension $d - 1$ that is orthogonal to $\boldsymbol{\mu}$. Therefore $\omega(\alpha, \theta, \varepsilon_0) = \omega_s(\mathcal{S}) = \min \left(\alpha, \sqrt{\frac{1}{\varepsilon_0^2 V^2} - \theta^2} \right)$ and the separability threshold reduces to

$$\delta_* = \max_{\alpha \geq 0, \theta \leq \frac{1}{\varepsilon_0 V}} \frac{\min \left(\alpha^2, \frac{1}{\varepsilon_0^2 V^2} - \theta^2 \right)}{\mathbb{E} \left[\left(1 - V\theta + \sqrt{\alpha^2 + \theta^2} g \right)_+^2 \right]}.$$

Note that the above ratio is decreasing in α over the range of $\alpha \geq \sqrt{\frac{1}{\varepsilon_0^2 \|\boldsymbol{\mu}\|_{\ell_2}^2} - \theta^2}$. Therefore, the maximizer α should satisfy $\alpha \leq \sqrt{\frac{1}{\varepsilon_0^2 \|\boldsymbol{\mu}\|_{\ell_2}^2} - \theta^2}$ and this further simplifies the expression for δ_* as follows

$$\delta_* = \max_{\alpha \geq 0, \alpha^2 + \theta^2 \leq \frac{1}{\varepsilon_0^2 V^2}} \frac{\alpha^2}{\mathbb{E} \left[\left(1 - V\theta + \sqrt{\alpha^2 + \theta^2} g \right)_+^2 \right]}.$$

By the change of variable $(\alpha, \frac{\theta}{\alpha}) \rightarrow (\alpha, \theta)$, this can be written as:

$$\delta_* = \left(\min_{\alpha \geq 0, \theta, \alpha^2(1+\theta^2) \leq \frac{1}{\varepsilon_0^2 V^2}} \mathbb{E} \left[\left(\frac{1}{\alpha} - V\theta + \sqrt{1 + \theta^2} g \right)_+^2 \right] \right)^{-1}.$$

Since the inner function is decreasing in α it is minimized at $\alpha_* = \frac{1}{\varepsilon_0 V \sqrt{1+\theta^2}}$ which simplifies the separability threshold to the following:

$$(4.4) \quad \delta_* = \left(\min_{\theta} \mathbb{E} \left[\left((\varepsilon_0 \sqrt{1+\theta^2} - \theta) V + \sqrt{1+\theta^2} g \right)_+^2 \right] \right)^{-1}.$$

To the best of our knowledge, our paper is the first work that shows such a phase transition for robust separability in the adversarial setting. In the non-adversarial case, similar phase transitions have been shown for data separability (a.k.a interpolation threshold) [8, 49, 12]. More specifically, [8] derived separability threshold for a logistic link regression model. Similar phenomenon extends to other link functions, as characterized by [49], and also to Gaussian mixtures model [12]. Interestingly, our result specialized to the case where the adversary has no power (cf. Example 1) recovers the existing thresholds for Gaussian mixtures model.

We end this section by demonstrating that in addition to the examples above Assumption 3 holds for a fairly broad family of mean vectors. This is the subject of the next lemma. We defer the proof of this lemma to Appendix E.2.

Assumption 4 *Suppose that the empirical distribution of the entries of $\sqrt{d}\boldsymbol{\mu}$ converges weakly to a distribution \mathbb{P}_M on real line, with bounded 2nd and p^{th} moment ($\int x^2 d\mathbb{P}_M(x) = \sigma_{M,2}^2 < \infty$, $\int |x|^p d\mathbb{P}_M(x) = \sigma_{M,p}^p < \infty$).*

Lemma 4.2 *Consider the asymptotic regime of $n \rightarrow \infty$ and $n/d \rightarrow \delta$, for some $\delta \in (0, \infty)$. Also, consider the function $J_q(\cdot; \cdot) : \mathbb{R} \times \mathbb{R}_{\geq 0} \mapsto \mathbb{R}_{\geq 0}$ defined by*

$$(4.5) \quad J_q(x; \lambda) = \min_u \frac{1}{2}(x - u)^2 + \lambda |u|^q.$$

Then Assumption 4 implies Assumption 3 with

$$(4.6) \quad \omega(\alpha, \theta, \varepsilon_0) = \min_{\lambda_0, \eta \geq 0, \nu} \sqrt{\delta} \left\{ \frac{\nu^2}{2\eta} + \frac{1}{2\eta\delta} + \frac{\eta}{2} \alpha^2 + \lambda_0 (\varepsilon_0 \sigma_{M,p})^{-q} \right\} \\ - \eta \sqrt{\delta} \mathbb{E} \left[J_q \left(\frac{h}{\eta \sqrt{\delta}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{M}{\sigma_{M,2}}; \frac{\lambda_0}{\eta} \right) \right],$$

where the expectation in the last line is taken with respect to the independent random variables $h \sim \mathcal{N}(0, 1)$ and $M \sim \mathbb{P}_M$.

4.2. Precise characterization of SA and RA in the separable regime In this section we precisely characterize the SA and RA of the classifier obtained as the limiting point of gradient descent on the loss (2.3) in the separable regime. As discussed in Proposition 3.2, the normalized iterations of gradient descent for the loss (2.3) converge to the max-margin classifier (F.2). Since $\text{SA}(\theta)$ and $\text{RA}(\theta)$ are only functions of the direction $\frac{\theta}{\|\theta\|_{\ell_2}}$, instead of studying the classifier obtained via GD iterations directly, we study the classification performance of the max-margin classifier.

Recall the function J_q is given by (4.5), and define

$$(4.7) \quad \mathcal{J}(c_0, c_1; \lambda_0) = \mathbb{E} \left[J_q \left(\frac{c_0}{\sqrt{\delta}} h - c_1 \frac{M}{\sigma_{M,2}}; \lambda_0 \sigma_{M,p}^q \right) \right],$$

where the expectation in the last line is taken with respect to the independent random variables $h \sim \mathcal{N}(0, 1)$ and $M \sim \mathbb{P}_M$, per the setting of Assumption 4. Our characterization of SA and RA will be in terms of the function \mathcal{J} as formalized in the next theorem.

Theorem 4.3 *Consider a data set generated i.i.d. according to an isotropic Gaussian mixture data model per Section 2.2 and suppose the mean vector μ obeys Assumptions 1 and 4. Also let $\tilde{\theta}^\varepsilon$ be the max margin solution per (F.2). If $\delta < \delta_*$, with δ_* given by (4.3), then in the asymptotic setting of Assumption 1 we have:*

- (a) *The following convex-concave minimax scalar optimization has a bounded solution $(\alpha_*, \gamma_{0*}, \theta_*, \beta_*, \lambda_{0*}, \eta_*, \tilde{\eta}_*)$ with the minimization components $(\alpha_*, \gamma_{0*}, \theta_*)$ unique:*

$$\min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\beta, \lambda_0, \eta \geq 0, \tilde{\eta}} D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta}), \quad \text{where}$$

$$\begin{aligned}
D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta}) &= 2 \left(1 + \frac{\eta}{2\alpha}\right)^{-1} \mathcal{J} \left(\frac{\beta}{2}, \frac{\tilde{\eta}}{2}; \frac{\lambda_0}{q\gamma_0^{q-1}} \left(1 + \frac{\eta}{2\alpha}\right)^{1-q} \right) \\
&\quad - \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right) \frac{1}{4(1 + \frac{\eta}{2\alpha})} - \frac{2\lambda_0}{q} \gamma_0 - \frac{\eta\alpha}{2} - \tilde{\eta}\theta \\
(4.8) \quad &\quad + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right]},
\end{aligned}$$

where the expectation in the last part is taken with respect to $g \sim \mathcal{N}(0, 1)$.

(b) It holds in probability that

$$(4.9) \quad \lim_{n \rightarrow \infty} \frac{1}{\|\boldsymbol{\mu}\|_{\ell_2}} \langle \boldsymbol{\mu}, \tilde{\boldsymbol{\theta}}^\varepsilon \rangle = \theta_*,$$

$$(4.10) \quad \lim_{n \rightarrow \infty} \|\tilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2} = \alpha_*,$$

$$(4.11) \quad \lim_{n \rightarrow \infty} \|\boldsymbol{\mu}\|_{\ell_p} \|\tilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_q} = \gamma_0 \alpha_*.$$

(c) Furthermore, part (b) combined with Lemma 2.1 imply the following limits hold in probability:

$$(4.12) \quad \lim_{n \rightarrow \infty} \text{SA}(\tilde{\boldsymbol{\theta}}^\varepsilon) = \Phi \left(\sigma_{M,2} \frac{\theta_*}{\alpha_*} \right),$$

$$(4.13) \quad \lim_{n \rightarrow \infty} \text{RA}(\tilde{\boldsymbol{\theta}}^\varepsilon) = \Phi \left(-\frac{\varepsilon_0 \gamma_0 \alpha_*}{\alpha_*} + \sigma_{M,2} \frac{\theta_*}{\alpha_*} \right).$$

Theorem 4.3 above provides us with a precise characterization of SA and RA and allows us to rigorously quantify the effect of adversary's manipulative power ε_0 , mean vector $\boldsymbol{\mu}$, and scaling of dimensions δ on SA and RA. In particular, this theorem precisely characterizes the performance of the max margin classifier (and in turn the classifier GD converges to) in terms of the optimal solutions to a low-dimensional optimization problem, namely (4.8). It is worth noting that by part (b), θ_* is the asymptotic value of the projection of the estimator $\tilde{\boldsymbol{\theta}}^\varepsilon$ along the direction of the class averages $\boldsymbol{\mu}$, and α_* represents the asymptotic value of the ℓ_2 norm of the estimator. Therefore, the θ_*/α_* term appearing in the SA and RA formulae corresponds to the correlation coefficient between the estimator $\tilde{\boldsymbol{\theta}}^\varepsilon$ and the class averages $\boldsymbol{\mu}$.

While the optimization problem (4.8) may look quite complicated, we note that it is a convex-concave problem in a handful number of scalar variables and hence can be easily solved by a low-dimensional gradient descent/ascent rather fast and accurately. In addition, this low-dimensional optimization problem significantly simplifies for special cases of p . We discuss some of these cases, which are also of particular practical interest, in Sections 5.1 and 5.2.

4.3. *Precise characterization of SA and RA in non-separable regime* In this section we precisely characterize the SA and RA of the classifier obtained by running gradient descent on the loss (2.3) in the non-separable regime. Before we can state our main result we need the definition of the Moreau envelop.

Definition 4.4 (Moreau envelope and expected Moreau envelope)

The Moreau envelope or Moreau-Yosida regularization of a function ℓ is given by

$$(4.14) \quad e_\ell(x; \mu) := \min_t \frac{1}{2\mu} (x - t)^2 + \ell(t).$$

We also define the expected Moreau envelope

$$(4.15) \quad L(a, b, \mu) = \mathbb{E}[e_\ell(ag + b; \mu)],$$

where the expectation is taken with respect to $g \sim \mathcal{N}(0, 1)$.

We this definition in place we are now ready to state our main result in the non-separable regime.

Theorem 4.5 Consider a data set generated i.i.d. according to an isotropic Gaussian mixture data model per Section 2.2 and suppose the mean vector μ obeys Assumption 4. Also let $\widehat{\theta}^\varepsilon$ be the solution to optimization (2.3). If $\delta > \delta_*$, with δ_* given by (4.3), then in the asymptotic setting of Assumption 1 we have:

- (a) The following convex-concave minimax scalar optimization has a bounded solution $(\theta_*, \alpha_*, \gamma_{0*}, \tau_{g*}, \beta_*, \tau_{h*})$ with the minimization components $(\alpha_*, \gamma_{0*}, \theta_*)$ unique:

$$(4.16) \quad \begin{aligned} & \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} \max_{0 \leq \beta, \tau_h} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) \\ & D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) = \frac{\beta \tau_g}{2} + L\left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2}\theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta}\right) \\ & - \min_{\lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^q + \frac{\nu^2}{2} - \mathcal{J}\left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \nu \right); \lambda_0\right) \right\} + \frac{\alpha \tau_h}{2} \right]. \end{aligned}$$

- (b) It holds in probability that

$$(4.17) \quad \lim_{n \rightarrow \infty} \frac{1}{\|\mu\|_{\ell_2}} \langle \mu, \widehat{\theta}^\varepsilon \rangle = \theta_*,$$

$$(4.18) \quad \lim_{n \rightarrow \infty} \|\mathbf{P}_\mu^\perp \widehat{\theta}^\varepsilon\|_{\ell_2} = \alpha_*,$$

$$(4.19) \quad \lim_{n \rightarrow \infty} \|\mu\|_{\ell_p} \|\widehat{\theta}^\varepsilon\|_{\ell_q} = \gamma_{0*}.$$

(c) As a corollary of part (b) and Lemma 2.1, the following limits hold in probability:

$$(4.20) \quad \lim_{n \rightarrow \infty} \text{SA}(\widehat{\boldsymbol{\theta}}^\varepsilon) = \Phi \left(\frac{\sigma_{M,2}\theta_*}{\sqrt{\alpha_*^2 + \theta_*^2}} \right),$$

$$(4.21) \quad \lim_{n \rightarrow \infty} \text{RA}(\widehat{\boldsymbol{\theta}}^\varepsilon) = \Phi \left(\frac{-\varepsilon_0 \gamma_{0*} + \sigma_{M,2}\theta_*}{\sqrt{\alpha_*^2 + \theta_*^2}} \right).$$

It is worth noting that by part (b), θ_* is the asymptotic value of the projection of the estimator $\widehat{\boldsymbol{\theta}}^\varepsilon$ along the direction of the class averages $\boldsymbol{\mu}$. In addition,

$$\lim_{n \rightarrow \infty} \|\mathbf{P}_{\boldsymbol{\mu}}^\perp \widehat{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2}^2 + \|\mathbf{P}_{\boldsymbol{\mu}} \widehat{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2}^2 = \alpha_*^2 + \theta_*^2$$

represents the asymptotic value of the squared ℓ_2 norm of the estimator. Therefore, the $\theta_*/\sqrt{\alpha_*^2 + \theta_*^2}$ term appearing in the SA and RA formulae corresponds to the correlation coefficient between the estimator $\widehat{\boldsymbol{\theta}}^\varepsilon$ and the class averages $\boldsymbol{\mu}$.

Theorem 4.5 complements the result of Theorem 4.3 by providing a precise characterization of SA and RA measures in the non-separable regime. In the remaining part of this section and also in the next section, we specialize our results to several specific choices of p that are of particular practical interest.

Remark 4.4 *In stating our results (Theorems 4.3 and 4.5), we are implicitly assuming the same variable ε_0 for both the perturbation level to the test data as well as the ‘perceived’ perturbation level used in the robust minimax estimator $\widehat{\boldsymbol{\theta}}^\varepsilon$. In principle, we can use different variable for the test perturbation level, say $\varepsilon_{0,\text{test}}$. The same results applies to this setting with minimal modifications; only in the RA formulae, cf. equations (4.13), (4.21) the variable ε_0 should be replaced by $\varepsilon_{0,\text{test}}$.*

5. Results for special cases of p In this section we discuss the implications of our main results for the special cases of ℓ_p perturbations with $p = 2$ in Section 5.1, $p = \infty$ in Section 5.2, and $p = 1$ in Section 5.3. We refer to Appendix D for the proofs of theorems and corollaries stated in this section.

5.1. Results for ℓ_2 perturbation We begin with stating our results for ℓ_p perturbation with $p = 2$. This result can be viewed as a corollary of Theorem 4.1, Theorem 4.3, and Theorem 4.5 specializing our main result for $p = 2$.

Corollary 5.1 *Consider a data set generated i.i.d. according to an isotropic Gaussian mixture data model per Section 2.2 and suppose the mean vector $\boldsymbol{\mu}$ obeys Assumptions 4. Then in the asymptotic setting of Assumption 1 we have:*

(a) The separability threshold δ_* is given by

$$(5.1) \quad \delta_* = \left(\min_{\theta} \mathbb{E} \left[\left((\varepsilon_0 \sqrt{1 + \theta^2} - \theta) V + \sqrt{1 + \theta^2} g \right)_+^2 \right] \right)^{-1}.$$

(b) In the separable regime where $\delta < \delta_*$, the followings hold in probability for the max margin solution $\tilde{\theta}^\varepsilon$ (see (F.2)):

$$(5.2) \quad \lim_{n \rightarrow \infty} \text{SA}(\tilde{\theta}^\varepsilon) = \Phi \left(\sigma_{M,2} \frac{\theta_*}{\alpha_*} \right), \quad \lim_{n \rightarrow \infty} \text{RA}(\tilde{\theta}^\varepsilon) = \Phi \left(\frac{\theta_*}{\alpha_*} \sigma_{M,2} - \varepsilon_0 \sigma_{M,2} \right),$$

where

$$(5.3) \quad \alpha_* = \left(\tilde{\alpha}_*^{-1} - \varepsilon_0 \sigma_{M,2} \right)^{-1}, \quad \theta_* = u_* \alpha_*.$$

Here, $(\tilde{\alpha}_*, u_*)$ the solution to the following problem:

$$(5.4) \quad \begin{aligned} & \min_{\tilde{\alpha} \geq 0, u} \tilde{\alpha}^2 \\ & \text{subject to} \quad 1 \geq u^2 + \delta \mathbb{E} \left[\left(\frac{1}{\tilde{\alpha}} - u \sigma_{M,2} + g \right)_+^2 \right], \end{aligned}$$

with expectation taken with respect to $g \sim \mathbf{N}(0, 1)$.

(c) In the non-separable regime where $\delta > \delta_*$, the followings hold in probability for the optimal solution $\hat{\theta}^\varepsilon$ of (2.3):

$$(5.5) \quad \lim_{n \rightarrow \infty} \text{SA}(\hat{\theta}^\varepsilon) = \Phi \left(\frac{\sigma_{M,2} \theta_*}{\sqrt{\alpha_*^2 + \theta_*^2}} \right),$$

$$(5.6) \quad \lim_{n \rightarrow \infty} \text{RA}(\hat{\theta}^\varepsilon) = \Phi \left(\frac{\theta_*}{\sqrt{\alpha_*^2 + \theta_*^2}} \sigma_{M,2} - \varepsilon_0 \sigma_{M,2} \right).$$

where $(\alpha_*, \theta_*, \beta_*)$ is the bounded solution of the following convex-concave minimax scalar optimization problem with the minimization components (α_*, θ_*) unique:

$$(5.7) \quad \begin{aligned} & \max_{0 \leq \beta} \min_{\theta, 0 \leq \alpha} D_{\text{ns}}(\alpha, \theta, \beta) \\ & D(\alpha, \theta, \beta) = L \left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2} \theta - \varepsilon_0 \sqrt{\alpha^2 + \theta^2}, \frac{\alpha}{\beta \sqrt{\delta}} \right) - \frac{\alpha \beta}{2 \sqrt{\delta}}. \end{aligned}$$

The corollary above precisely characterizes the behavior of the classifier that gradient descent converges to in terms of low-dimensional optimization problems ((5.4) in the separable regime and (5.7) in the non-separable regime).

Recall that the term θ_*/α_* in the separable regime and the term $\theta_*/\sqrt{\alpha_*^2 + \theta_*^2}$ in the non-separable regime correspond to the correlation coefficient between the robust minimax estimator and the classes average $\boldsymbol{\mu}$. As we will see in Figure 4, the standard accuracy is decreasing in $1/\delta$, for any fixed ε_0 , which equivalently indicates that the correlation between the estimator and $\boldsymbol{\mu}$ is monotone increasing in the sample-to-dimension ratio δ .

As we will see in the coming sections, SA and RA curves have a highly non-trivial behavior which also strongly depend on the choice of p . This necessitate a rigorous theory (such as the above) that can precisely predict these curves. To better understand the implications and consequences of this result we focus on its various predictions. Specifically, we find the global optima of the two low-dimensional optimization problems via simple gradient descent/ascent and use it to calculate the corresponding SA and RA based on (5.2) and (5.5). We also verify these theoretical predictions with the performance of gradient descent on the loss (2.3) with a polyak/approximate polyak step size in the separable/non-separable regimes.³

We plot the theoretically predicted standard and robust accuracy versus the adversary's power ε_0 together with the corresponding empirical results in Figure 3 (a) and (b). The solid lines depict theoretical predictions with the dots representing the empirical performance of gradient descent with the algorithmic settings discussed above. The data set is generated according to a Gaussian Mixture Model per Section 2.2 with $\boldsymbol{\mu} \in \mathbb{R}^d$ consisting of i.i.d. $\mathcal{N}(0, \frac{1}{d})$ entries with dimension $d = 400$. Each dot represents the average of 100 trials. These figures demonstrate that even for moderate dimension sizes our theoretical prediction is a near perfect match with the empirical performance of gradient descent. We note that when ε_0 is sufficiently large then the adversarially trained model $\hat{\boldsymbol{\theta}}^\varepsilon$ becomes zero due to the large regularization in the argument of loss function in (2.3) and SA and RA measures are not defined. The curves are plotted up to that ε_0 .

An intriguing observation of Corollary 5.1 is that in the separable regime in the case of $p = 2$, the standard accuracy does not depend on ε_0 . In other words, adversarial training has no effect on the performance on benign unperturbed data. The robust accuracy, however is decreasing in ε_0 . Figure 3 (a) and (b) also verify this predicted behavior and capture the effect of the adversary's power ε_0 on standard and robust accuracy. In the separable regime, SA is flat which implies that adversarial training has no effect on

³Specifically we run gradient descent iterations of the form $\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_\tau - \alpha_\tau \nabla \mathcal{L}(\boldsymbol{\theta}_\tau)$ on (2.3) with a Polyak step size $\alpha_\tau = \frac{\mathcal{L}(\boldsymbol{\theta}_\tau)}{\|\nabla \mathcal{L}(\boldsymbol{\theta}_\tau)\|_{\ell_2}^2}$ in the separable regime and an approximate Polyak step size $\alpha_\tau = \frac{\mathcal{L}(\boldsymbol{\theta}_\tau) - \min_{0 \leq t \leq \tau} \mathcal{L}(\boldsymbol{\theta}_t) + \frac{\gamma}{\tau}}{\|\nabla \mathcal{L}(\boldsymbol{\theta}_\tau)\|_{\ell_2}^2}$ in the non-separable regime.

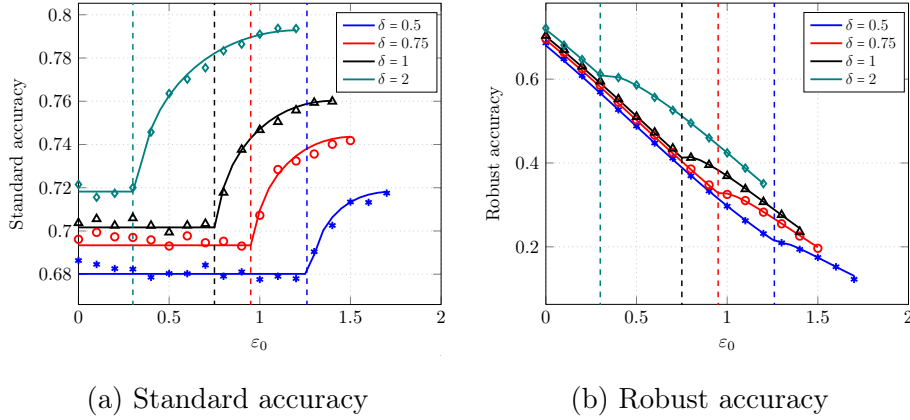


Fig 3: Depiction of standard and robust accuracies as a function of the adversary’s normalized power ε_0 with ℓ_2 ($p = 2$) perturbation for different values of δ . Solid curves are theoretical predictions and dots are results obtained based on gradient descent on the robust objective (2.3). The dashed lines depict the separability threshold for that δ . Each dot represents the average of 100 trials. The data set is generated according to a Gaussian Mixture Model per Section 2.2 with $\boldsymbol{\mu} \in \mathbb{R}^d$ consisting of i.i.d. $\mathcal{N}(0, \frac{1}{d})$ entries with dimension $d = 400$.

standard accuracy (or the generalization error on unperturbed data). However, adversarial training does affect RA because now the trained model is used to classify the adversarially perturbed test data.

In the non-separable regime, we observe that adversarial training helps with improving the standard accuracy! Further, such positive impact is observed for all choice of δ with a rather robust trend. Note that this behavior is significantly different from a regression setting where adversarial training first improves with the standard accuracy but then there is a turning point beyond which the standard accuracy will decrease as ε_0 grows. We refer to [30, Figure 3] and discussion therein for more details on a regression setting. Moreover, as depicted in Figure 3(b) we see that RA always declines as adversary gets more powerful (i.e., ε_0 grows) as expected.

Next in Figure 4, we plot SA and RA versus dimension-to-sample ratio $\frac{1}{\delta} = \frac{d}{n}$, which is a measure of model complexity, for several values of ε_0 . It has been shown that the standard risk (which amounts to $1 - \text{SA}$ in our setting) as a function of model complexity $\frac{1}{\delta}$ undergoes a double-descent behavior for various learning models [6, 5, 24]. Specifically, the risk depicts a

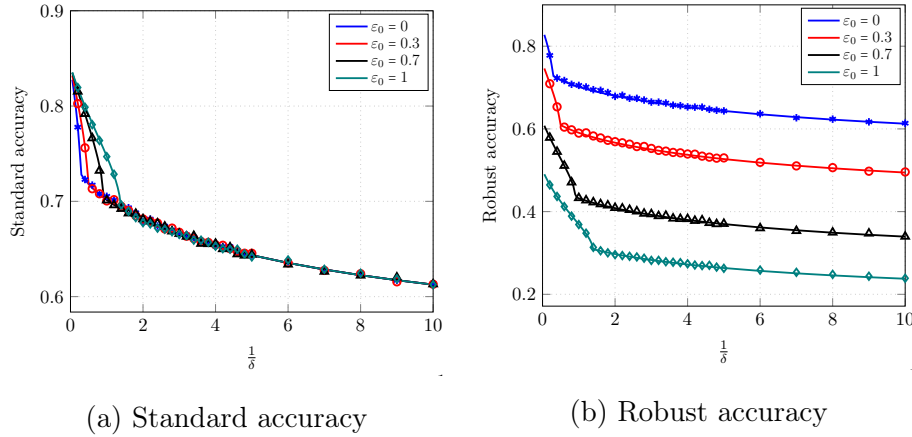


Fig 4: Depiction of standard and robust accuracies as a function of dimension-to-sample ratio $\frac{1}{\delta} = \frac{d}{n}$, which is a measure of model complexity, for several values of ε_0 with ℓ_2 ($p = 2$) perturbation, under a similar setting as in Figure 3.

U-shape before the interpolation threshold (separability threshold in binary classification) and then starts to decline afterwards. Interestingly, for the current setting of experiments here we do not observe such double descent behavior and the standard accuracy always decreases as $\frac{1}{\delta}$ grows, albeit at different rates in the separable and non-separable regimes.⁴

5.2. Results for ℓ_∞ perturbation For the case of ℓ_∞ perturbation ($p = \infty$ and $q = 1$), Theorem 4.1, Theorem 4.3, and Theorem 4.5 do not substantially simplify. However, we can calculate the function J_q defined by (4.5) in closed form. In this case J_q becomes the Huber function given by

$$J_1(x, \lambda) = \begin{cases} \lambda|x| - \frac{\lambda^2}{2} & |x| \geq \lambda \\ \frac{x^2}{2} & |x| \leq \lambda \end{cases}$$

Using Theorems 4.1, 4.3, and 4.5 with this closed form for J_q , in Figure 5, we again depict our theoretical predictions for standard and robust accuracy as well as the empirical performance of gradient descent as a function of the

⁴It is worth noting that the double descent phenomenon has been observed for binary classification in a non-adversarial setting with model misspecification. In such a model the learner observes only a subset $S \subset [d]$ of size p of the covariates with $d/n \rightarrow \zeta \geq 1$ and $p/n \rightarrow \kappa \in (0, \zeta]$ (see [12] for further details). Our theoretical analysis can in principle be used to analyze such a setting, however we do not pursue this direction in this paper.

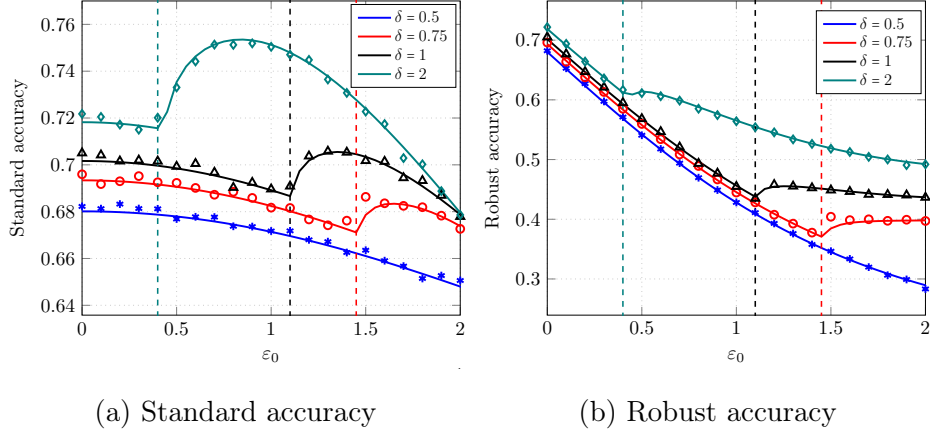


Fig 5: Depiction of standard and robust accuracies as a function of ε_0 with ℓ_∞ ($p = \infty$) perturbation for different values of δ , and under a similar setting as in Figure 3.

adversary's normalized power for various values of δ . As in the $p = 2$ case our theoretical predictions is very accurate even for moderate dimensions d .

More specifically, Figure 5(a) depicts the standard accuracy (SA) versus the adversary's normalized power. Similar to our $p = 2$ results the data set is generated according to a Gaussian Mixture Model per Section 2.2 with $\mu \in \mathbb{R}^d$ consisting of i.i.d. $\mathcal{N}(0, \frac{1}{d})$ entries with dimension $d = 400$ and each data points represents the average of 100 trials. In the case of $p = \infty$ however, we do not use the scaling $\varepsilon = \varepsilon_0 \|\mu\|_{\ell_\infty}$ as $\|\sqrt{d}\mu\|_{\ell_\infty}$ grows with $\sqrt{\log d}$ and therefore violates Assumption 4. Instead we shall use a slightly different scaling of $\varepsilon = \frac{\varepsilon_0}{\sqrt{d}}$. In the separable regime, we see that adversarial training hurts the standard accuracy. However, in the non-separable regime, the standard accuracy starts increasing indicating that adversarial training is improving the standard accuracy. Furthermore, after some value of ε_0 , which interestingly shifts with δ , the standard accuracy starts to go down as ε_0 grows.⁵ We note that this behavior is rather counterintuitive and very different from the $p = 2$ case, further highlighting the need for a precise theory that can predict such nuanced behavior. Figure 5(b) shows the robust accuracy RA versus ε_0 for various values of δ . In the separable regime, we observe a similar trend for all δ , namely RA decreases at an almost linear rate. In the non-separable regime though we have different trends depending on the value of δ .

⁵Note that for $\delta = 0.5$, we are in the separable regime over the entire range $[0, \varepsilon_0]$.

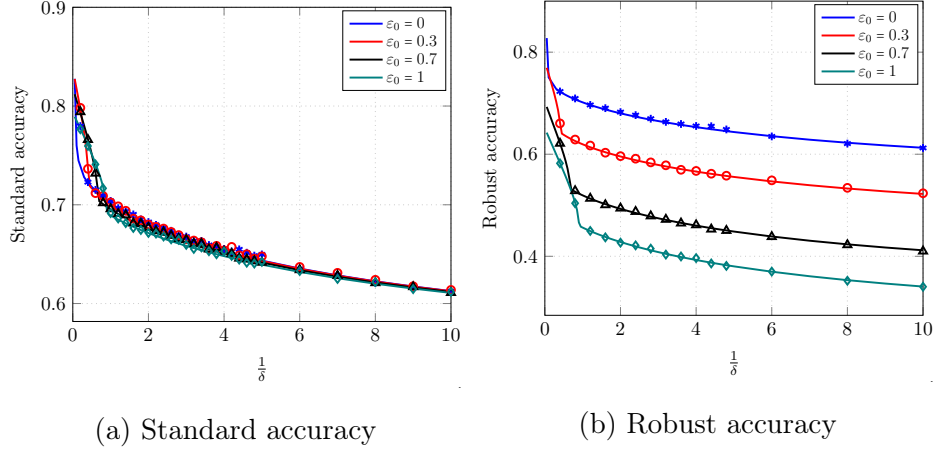


Fig 6: Depiction of standard and robust accuracies as a function of dimension-to-sample ratio $\frac{1}{\delta} = \frac{d}{n}$, which is a measure of model complexity, for several values of ε_0 with ℓ_∞ ($p = \infty$) perturbation, under a similar setting as in Figure 3.

Finally, in Figure 6 we depict the effect of overparameterization $\frac{1}{\delta}$ on SA and RA. We observe a similar pattern as in the case of $p = 2$. In particular, we do not observe a double descent behavior and the standard accuracy always decreases as $\frac{1}{\delta}$ grows, albeit at different rates in the separable and non-separable regimes.

5.3. Results for ℓ_1 perturbation Our characterization of SA and RA given by Theorem 4.3, for separable regime, and by Theorem 4.5, for non-separable regime involve the function \mathcal{J} defined by (4.7) which in turn depends on the function J_q given by (4.5). However, J_q is only defined for finite q and therefore the case of $p = 1$, $q = \infty$ is not directly covered by our results in Section 4. That said, a very similar analysis can be used to characterize SA and RA in this case. We formalize our results for this case in the next theorem.

Theorem 5.2 *Consider a data set generated i.i.d. according to an isotropic Gaussian mixture data model per Section 2.2 and suppose the mean vector μ obeys Assumptions 4. Also define*

$$(5.8) \quad f(c_0, c_1; t_0) = \frac{1}{2} \mathbb{E} \left[\text{ST} \left(\frac{c_0}{\sqrt{\delta}} h - c_1 \frac{M}{\sigma_{M,2}}; \frac{t_0}{\sigma_{M,1}} \right)^2 \right],$$

where $\text{ST}(x; a) := \text{sgn}(x) (|x| - a)_+$ is the soft-thresholding function. Then in the asymptotic setting of Assumption 1 we have:

(a) The separability threshold δ_* is given by

$$\delta_* := \min_{\alpha \geq 0, \theta} \frac{\omega(\alpha, \theta, \varepsilon_0)^2}{\mathbb{E} \left[\left(1 - V\theta + \sqrt{\alpha^2 + \theta^2} g \right)_+^2 \right]} \quad (5.9)$$

with $\omega(\alpha, \theta, \varepsilon_0) := \min_{\eta \geq 0, \nu} \sqrt{\delta} \left\{ \frac{\nu^2}{2\eta} + \frac{1}{2\eta\delta} + \frac{\eta}{2} \alpha^2 - \eta f \left(\frac{1}{\eta}, \frac{\nu}{\eta} - \theta; \frac{1}{\varepsilon_0} \right) \right\}.$

(b) In the separable regime where $\delta < \delta_*$, the followings hold in probability for the max margin solution $\tilde{\theta}^\varepsilon$ (see (F.2)):

$$(5.10) \quad \lim_{n \rightarrow \infty} \text{SA}(\tilde{\theta}^\varepsilon) = \Phi \left(\sigma_{M,2} \frac{\theta_*}{\alpha_*} \right),$$

$$(5.11) \quad \lim_{n \rightarrow \infty} \text{RA}(\tilde{\theta}^\varepsilon) = \Phi \left(-\frac{\varepsilon_0 \gamma_{0*}}{\alpha_*} + \sigma_{M,2} \frac{\theta_*}{\alpha_*} \right).$$

Here, $(\alpha_*, \gamma_{0*}, \theta_*)$ are the unique minimization component of the following convex-concave minimax scalar optimization with bounded solution $(\alpha_*, \gamma_{0*}, \theta_*, \beta_*, \eta_*, \tilde{\eta}_*)$.

$$\min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\beta, \eta \geq 0, \tilde{\eta}} D_s(\alpha, \gamma_0, \theta, \beta, \eta, \tilde{\eta}), \quad \text{where}$$

$$D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta}) = \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\beta, \eta \geq 0, \tilde{\eta}} \frac{1}{2(1 + \frac{\eta}{2\alpha})} f \left(\beta, \tilde{\eta}; 2\gamma_0 \left(1 + \frac{\eta}{2\alpha} \right) \right) - \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right) \frac{1}{4(1 + \frac{\eta}{2\alpha})} - \frac{\eta\alpha}{2} - \tilde{\eta}\theta$$

$$(5.12) \quad + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right]},$$

with expectation taken with respect to $g \sim \mathcal{N}(0, 1)$.

(c) In the non-separable regime where $\delta > \delta_*$, the followings hold in probability the optimal solution $\tilde{\theta}^\varepsilon$ of (2.3):

$$(5.13) \quad \lim_{n \rightarrow \infty} \text{SA}(\tilde{\theta}^\varepsilon) = \Phi \left(\frac{\sigma_{M,2} \theta_*}{\sqrt{\alpha_*^2 + \theta_*^2}} \right),$$

$$(5.14) \quad \lim_{n \rightarrow \infty} \text{RA}(\tilde{\theta}^\varepsilon) = \Phi \left(\frac{-\varepsilon_0 \gamma_{0*} + \sigma_{M,2} \theta_*}{\sqrt{\alpha_*^2 + \theta_*^2}} \right).$$

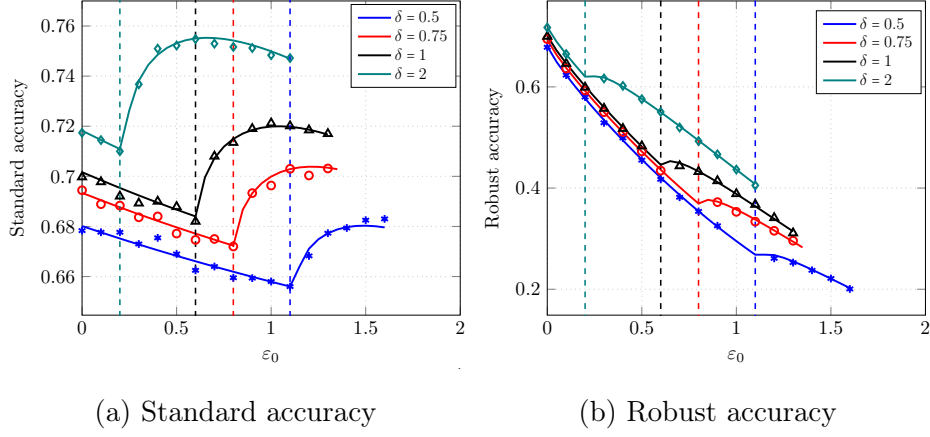


Fig 7: Depiction of standard and robust accuracies as a function of ε_0 with ℓ_1 ($p = 1$) perturbation for different values of δ , under a similar setting as in Figure 3.

Here, $(\alpha_*, \gamma_{0*}, \theta_*)$ are the unique minimization components of the following convex-concave minimax scalar optimization with bounded solution $(\theta_*, \alpha_*, \gamma_{0*}, \tau_{g*}, \beta_*, \tau_{h*})$.

$$\begin{aligned}
 & \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} \max_{0 \leq \beta, \tau_h} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) \\
 D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) &= \frac{\beta \tau_g}{2} + L \left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2} \theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta} \right) \\
 (5.15) \quad & - \min_{\nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \frac{\nu^2}{2} - f \left(\beta, \frac{\tau_h \theta}{\alpha} + \nu; \frac{\gamma_0 \tau_h}{\alpha} \right) \right\} + \frac{\alpha \tau_h}{2} \right].
 \end{aligned}$$

In Figure 7, we again depict our theoretical predictions for standard and robust accuracy as well as the empirical performance of gradient descent as a function of the adversary's normalized power for various values of δ . We note however that in this case we do not actually run gradient descent in our simulations as $p = 1$ corresponds to $q = +\infty$ and GD convergence is extremely slow since the gradient only has one non-zero entry. Therefore, for our empirical simulations we use CVX, a package for specifying and solving convex programs [22], in the non-separable regime which given the uniqueness of the global optima yields the same answer as GD. Similarly, in the separable regime we use (F.2) which based on Proposition 3.2 is the direction GD eventually converges to. We observe that as in the $p = 2$ and $p = +\infty$ cases our theoretical predictions are very accurate even for moderate dimensions d .

More specifically, Figure 7(a) depicts the standard accuracy (SA) versus the adversary's normalized power. Similar to our $p = +\infty$ results the data set is generated according to a Gaussian Mixture Model per Section 2.2 with $\boldsymbol{\mu} \in \mathbb{R}^d$ consisting of i.i.d. $\mathcal{N}(0, \frac{1}{d})$ entries with dimension $d = 400$ and each data points represents the average of 100 trials. In the separable regime, we see that adversarial training hurts the standard accuracy. However, in the non-separable regime, the standard accuracy starts increasing indicating that adversarial training is improving the standard accuracy. Furthermore, after some value of ε_0 , which interestingly shifts with δ , the standard accuracy starts to go down as ε_0 grows.⁶ We note that this behavior is rather counterintuitive and very different from the $p = 2$ case but somewhat similar to the $p = +\infty$ case. This again highlights the need for a precise theory that can predict such nuanced behavior. Figure 7(b) shows the robust accuracy RA versus ε_0 for various values of δ . In the separable regime, we observe a similar trend for all δ , namely RA decreases at an almost linear rate. In the non-separable regime though we have different trends depending on the value of δ .

6. Extension to anisotropic Gaussians In this section we extend our results to Gaussian distributions with general covariance matrices that obey a certain spiked covariance assumption stated below.

Assumption 5 (*Spiked covariance*) $\boldsymbol{\mu}$ is an eigenvector of $\boldsymbol{\Sigma}$ with eigenvalue a^2 , i.e., $\boldsymbol{\Sigma}\boldsymbol{\mu} = a^2\boldsymbol{\mu}$.

Similar spiked covariance models have been used to model data in a number of statistical problems, including matrix denoising and structured learning [32, 16], sparse PCA [13], synchronization and clustering [28].

To extend our results in Section 4 to the anisotropic case we also need to generalize the definition of the set \mathcal{S} as follows:

$$\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu}) := \left\{ \mathbf{z} \in \mathbb{R}^d : \quad \mathbf{z}^T \tilde{\boldsymbol{\mu}} = 0, \|\mathbf{z}\|_{\ell_2} = \alpha, \left\| \boldsymbol{\Sigma}^{-1/2} \mathbf{z} + \theta \tilde{\boldsymbol{\mu}} \right\|_{\ell_q} \leq \frac{1}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}} \right\}.$$

We are now ready to state our main results in the anisotropic case. We start by the separability threshold which generalizes Theorem 4.1.

Theorem 6.1 *Consider a data set generated i.i.d. according to an anisotropic Gaussian mixture data model per Section 2.2 with a spiked covariance per Assumption 5. Also suppose the mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ obey*

⁶Note that for $\delta = 0.5$, we are in the separable regime over the entire range $[0, \varepsilon_0]$.

Assumptions 2 and 3. Also define

$$(6.1) \quad \delta_* := \min_{\alpha \geq 0, \theta} \frac{\omega(\alpha, \theta, \varepsilon_0)^2}{\mathbb{E} \left[\left(1 - V\theta + \sqrt{\alpha^2 + a^2 \theta^2} g \right)_+^2 \right]},$$

where the expectation is taken with respect to $g \sim \mathcal{N}(0, 1)$. Then, under the asymptotic setting of Assumption 1, for $\delta < \delta_*$ the data are (ε, q) -separable with high probability and for $\delta > \delta_*$, the data are non-separable, with high probability. Namely,

$$\begin{aligned} \delta < \delta_* &\Rightarrow \lim_{n \rightarrow \infty} \mathbb{P}(\text{data is } (\varepsilon, q)\text{-separable}) = 1, \\ \delta > \delta_* &\Rightarrow \lim_{n \rightarrow \infty} \mathbb{P}(\text{data is } (\varepsilon, q)\text{-separable}) = 0. \end{aligned}$$

Our next theorem precisely characterizes SA and RA in the separable regime and generalizes Theorem 4.3 to the anisotropic case. Before proceeding to state the theorem we need to establish some definitions and assumptions.

Definition 6.2 For a given matrix $\mathbf{A} \geq 0$ and a function f , we define the weighted Moreau envelope of f as follows:

$$e_{f, \mathbf{A}}(\mathbf{x}; \lambda) := \min_{\mathbf{v}} \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_{\mathbf{A}}^2 + \lambda f(\mathbf{v})$$

When $\mathbf{A} = \mathbf{I}$, we recover the (scaled) classical Moreau envelope. We denote by $e_{q, \Sigma}$ the weighted Moreau envelope corresponding to $\|\cdot\|_{\ell_q}^q$ function.

Assumption 6 For the sequence of instances $\{\Sigma(n), \mu(n), d(n)\}_{n \in \mathbb{N}}$ indexed by n , we assume that:

- (a) The following (in probability) limit exists for any scalars $c_0, c_1, \lambda_0, \eta \in \mathbb{R}_+$:

$$F(c_0, c_1; b_0, b_1) := \lim_{n \rightarrow \infty} e_{q, \mathbf{I} + b_0 \Sigma} \left((\mathbf{I} + b_0 \Sigma)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \tilde{\mu} \right\}; b_1 \|\mu\|_{\ell_p}^q \right).$$

- (b) The empirical distribution of eigenvalues of Σ converges weakly to a distribution ρ with Stieltjes transform $S_{\rho}(z) := \int \frac{\rho(t)}{z-t} dt$.

With these definitions and assumptions in place we are ready to state our result in the separable regime.

Theorem 6.3 Consider a data set generated i.i.d. according to an anisotropic Gaussian mixture data model per Section 2.2 with a spiked covariance per Assumption 5. Also suppose the mean vector μ and covariance matrix Σ obey

Assumptions 2, 3, and 6. Also let $\tilde{\boldsymbol{\theta}}^\varepsilon$ be the max margin solution per (F.2). If $\delta < \delta_*$, with δ_* given by (4.3), then in the asymptotic setting of Assumption 1 we have:

- (a) The following convex-concave minimax scalar optimization problem has bounded solution $(\alpha_*, \gamma_{0*}, \theta_*, \beta_*, \lambda_{0*}, \eta_*, \tilde{\eta}_*)$ with the minimization components $(\alpha_*, \gamma_{0*}, \theta_*)$ unique:

$$\begin{aligned} \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\beta, \lambda_0, \eta \geq 0, \tilde{\eta}} D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta}), \quad \text{where} \\ D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta}) = 2F\left(\beta, \tilde{\eta}; \frac{\eta}{2\alpha}, \frac{\lambda_0}{q\gamma_0^{q-1}}\right) - \frac{\beta^2\alpha}{2\delta\eta} \left(1 + \frac{2\alpha}{\eta} S_\rho\left(-\frac{2\alpha}{\eta}\right)\right) \\ - \frac{2\lambda_0}{q}\gamma_0 - \frac{\eta\alpha}{2} - \tilde{\eta}\theta \\ - \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\alpha}a^2)} + \beta \sqrt{\mathbb{E}\left[\left((1 + \varepsilon_0\gamma_0 - \theta V) + \alpha g\right)_+^2\right]}, \end{aligned} \quad (6.2)$$

with expectation in last part taken with respect to $g \sim \mathbf{N}(0, 1)$.

- (b) It holds in probability that

$$(6.3) \quad \lim_{n \rightarrow \infty} \frac{1}{\|\boldsymbol{\mu}\|_{\ell_2}} \langle \boldsymbol{\mu}, \tilde{\boldsymbol{\theta}}^\varepsilon \rangle = \theta_*,$$

$$(6.4) \quad \lim_{n \rightarrow \infty} \|\tilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2} = \alpha_*,$$

$$(6.5) \quad \lim_{n \rightarrow \infty} \|\boldsymbol{\mu}\|_{\ell_p} \|\tilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_q} = \gamma_{0*}.$$

- (c) As a corollary of part (b) and Lemma 2.1, the following limits hold in probability:

$$(6.6) \quad \lim_{n \rightarrow \infty} \text{SA}(\tilde{\boldsymbol{\theta}}^\varepsilon) = \Phi\left(V \frac{\theta_*}{\alpha_*}\right),$$

$$(6.7) \quad \lim_{n \rightarrow \infty} \text{RA}(\tilde{\boldsymbol{\theta}}^\varepsilon) = \Phi\left(-\frac{\varepsilon_0\gamma_{0*}}{\alpha_*} + V \frac{\theta_*}{\alpha_*}\right).$$

Next we turn our attention to characterizing SA and RA on the non-separable regime. To state result we need an additional assumption

Assumption 7 For the sequence of instances $\{\boldsymbol{\Sigma}(n), \boldsymbol{\mu}(n), p(n)\}_{n \in \mathbb{N}}$ indexed by n , we assume that the following (in probability) limit exists for any scalars $c_0, c_1 \in \mathbb{R}_+$ and $\lambda_0 \in \mathbb{R}$:

$$(6.8) \quad \mathbb{E}(c_0, c_1; \lambda_0) := \lim_{n \rightarrow \infty} e_{q, \boldsymbol{\Sigma}} \left(\frac{c_0}{\sqrt{n}} \boldsymbol{\Sigma}^{-1/2} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_p}^q \right),$$

where we recall $\tilde{\boldsymbol{\mu}} = \boldsymbol{\mu} / \|\boldsymbol{\mu}\|_{\ell_2}$.

Our next theorem generalizes Theorem 4.5 to anisotropic case.

Theorem 6.4 *Consider a data set generated i.i.d. according to an anisotropic Gaussian mixture data model per Section 2.2 with a spiked covariance per Assumption 5. Also suppose the mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ obey Assumptions 2, 3, and 7. Also let $\widehat{\boldsymbol{\theta}}^\varepsilon$ be the solution to optimization (2.3). If $\delta > \delta_*$, with δ_* given by (4.3), then in the asymptotic setting of Assumption 1 we have:*

- (a) *The following convex-concave minimax scalar optimization problem has bounded solution $(\theta_*, \alpha_*, \gamma_{0*}, \tau_{g*}, \beta_*, \tau_{h*})$ with the minimization components $(\alpha_*, \gamma_{0*}, \theta_*)$ unique:*

$$\begin{aligned} & \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} \max_{0 \leq \beta, \tau_h} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) \\ D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) &= \frac{\beta \tau_g}{2} + L \left(\sqrt{\alpha^2 + a^2 \theta^2}, V\theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta} \right) \\ (6.9) \quad & - \min_{\lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^q + \frac{\nu^2}{2} - \mathbb{E} \left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right); \lambda_0 \right) \right\} + \frac{\alpha \tau_h}{2} \right]. \end{aligned}$$

- (b) *It holds in probability that*

$$(6.10) \quad \lim_{n \rightarrow \infty} \frac{1}{\|\boldsymbol{\mu}\|_{\ell_2}} \langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}}^\varepsilon \rangle = \theta_*,$$

$$(6.11) \quad \lim_{n \rightarrow \infty} \|\mathbf{P}_{\boldsymbol{\mu}}^\perp \widehat{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2} = \alpha_*,$$

$$(6.12) \quad \lim_{n \rightarrow \infty} \|\boldsymbol{\mu}\|_{\ell_p} \|\widehat{\boldsymbol{\theta}}^\varepsilon\|_{\ell_q} = \gamma_{0*}.$$

- (c) *As a corollary of part (b) and Lemma 2.1, the following limits hold in probability:*

$$(6.13) \quad \lim_{n \rightarrow \infty} \text{SA}(\widehat{\boldsymbol{\theta}}^\varepsilon) = \Phi \left(\frac{V\theta_*}{\sqrt{\alpha_*^2 + a^2 \theta_*^2}} \right),$$

$$(6.14) \quad \lim_{n \rightarrow \infty} \text{RA}(\widehat{\boldsymbol{\theta}}^\varepsilon) = \Phi \left(\frac{-\varepsilon_0 \gamma_{0*} + V\theta_*}{\sqrt{\alpha_*^2 + a^2 \theta_*^2}} \right).$$

The results above generalize out results to the anisotropic case. The reader may of course be wondering when Assumptions 6 and 7 hold. This is the subject of the next Remark which we prove in Appendix B.4.

Remark 6.1 *For the case of ℓ_2 perturbation ($p = q = 2$), the following two conditions are sufficient for Assumption 6 and 7 to hold:*

- (i) *The empirical distribution of the entries of $\sqrt{d}\boldsymbol{\mu}$ converges weakly to a distribution \mathbb{P}_M on real line, with bounded second moment, i.e. $\int x^2 d\mathbb{P}_M(x) = \sigma_{M,2}^2 < \infty$.*
- (ii) *The empirical distribution of eigenvalues of $\boldsymbol{\Sigma}$ converges weakly to a distribution ρ with Stieltjes transform $S_\rho(z) := \int \frac{\rho(t)}{z-t} dt$.*

7. Proof sketch and mathematical challenges Our theoretical results on adversarial training for binary classification fits in the rapidly growing recent literature on developing *sharp* high-dimensional asymptotics of (possibly non-smooth) convex optimization-based estimators [18, 59, 4, 1, 60, 17, 51, 64, 34, 19, 15, 50, 48, 67, 9, 25, 7]. Most of this line of work focus on linear models and regression problems. It has been only recently that the literature witnessed a surge of interest in sharp analysis of a variety of methods tailored to binary classification models [26, 8, 61, 43, 33, 56, 62, 12, 49, 38, 46, 62, 40]. However, none of these papers study adversarial training and its impact on standard/robust accuracies.

On a technical level, our sharp analysis relies on the Convex Gaussian Min-max Theorem (CGMT) [64] (see also [60, 51, 50]), which is a powerful extension of the Gordon’s Gaussian comparison inequality [21]. We refer to Section 7 for an overview of this framework and the mathematical challenges we encounter in applying it to our adversarial setting. We next present a proof sketch for deriving our main results which illustrates the key ideas.

To be able to provide a precise characterization of the various tradeoffs we need to develop a precise understanding of the adversarial training objective

$$(7.1) \quad \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q}),$$

and its optimal solution $\widehat{\boldsymbol{\theta}}^\varepsilon \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta})$. Given the classification nature of the problem, as discussed earlier, we have to study this loss in the two different regimes of separable and non-separable as well as characterize the threshold of separability. In this section we wish to provide a brief overview of the steps of our proofs and some of the challenges. We focus our exposition on the non-separable case. While the details of the derivations for the separable case and the calculation of the separability threshold differ from the non-separable case the general steps are similar and therefore the steps below also provides a general road map for the proof of these results as well. Specifically, our proofs in the non-separable regime consists of the following steps:

Step I: Reformulation of the loss.

The loss (7.1), while significantly simplified due to the removal of the max function, is still rather complicated and precisely characterizing the behavior

and the quality of its optimal solution is still challenging. In particular, the dependence on the random data matrix \mathbf{X} is still rather complex hindering statistical analysis even in an asymptotic setting. To bring the optimization problem into a form more amenable to precise asymptotic analysis we carry out a series of reformulations of the optimization problem. Combining these reformulation steps we arrive at the following equivalent Primal Optimization (PO) problem

$$(7.2) \quad \min_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^n} \max_{\mathbf{u} \in \mathbb{R}^n} \frac{1}{n} \left\{ \mathbf{u}^\top \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} + \mathbf{u}^\top \mathbf{D}_y \mathbf{Z} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} - \mathbf{u}^\top \mathbf{v} \right\} + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right)$$

Step II: Reduction to an Auxiliary Optimization (AO) problem.

The equivalent form above may be counter-intuitive as we started by simplifying a different mini-max optimization problem and we have now again introduced a new maximization! The main advantage of this new form is that it is in fact affine in the data matrix \mathbf{X} . This particular form allows us to use a powerful extension of a classical Gaussian process inequality due to [21] known as Convex Gaussian Minimax Theorem (CGMT) [64] which focuses on characterizing the asymptotic behavior of mini-max optimization problems that are affine in a Gaussian matrix \mathbf{X} . This result enables us to characterize the properties of (7.1) by studying the asymptotic behavior of the following, arguable simpler, *Auxiliary Optimization (AO)* problem instead

$$(7.3) \quad \min_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^n} \max_{\mathbf{u} \in \mathbb{R}^n} \frac{1}{n} \left\{ \left\| \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g}^T \mathbf{D}_y \mathbf{u} + \left\| \mathbf{D}_y \mathbf{u} \right\|_{\ell_2} \mathbf{h}^T \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right. \\ \left. + (\mathbf{u}^\top \mathbf{D}_y \mathbf{z}) (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) + \mathbf{u}^\top \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} - \mathbf{u}^\top \mathbf{v} \right\} + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right),$$

where $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_n)$ and $\mathbf{h} \sim \mathcal{N}(0, \mathbf{I}_d)$, $\mathbf{P}_\mu^\perp := \mathbf{I} - \tilde{\boldsymbol{\mu}} \tilde{\boldsymbol{\mu}}^T$, and $\mathbf{P}_\mu := \tilde{\boldsymbol{\mu}} \tilde{\boldsymbol{\mu}}^T$.

We emphasize that the relationship between the above PO problem (7.2) and how it is exactly related to the AO problem (7.3) is more intricate and technical compared with classical CGMT and related work in the context of classification [56, 62]. In particular, prior work on binary classification such as [56, 62] via CGMT (which corresponds to the non-robust case i.e. $\varepsilon = 0$) utilize the fact that (7.2) is rotationally invariant and hence one can assume $\boldsymbol{\mu} = \mathbf{e}_1$ without loss of generality. However, in the robust version (unless $p = q = 2$) the direction of $\boldsymbol{\mu}$ plays a crucial role due to the regularization term $\frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right)$.

Step III: Scalarization of the Auxiliary Optimization (AO) problem.

In this step we further simplify the AO problem in (7.3). In particular we

show the asymptotic behavior of the AO can be characterized rather precisely via the scalar optimization problem

$$(7.4) \quad \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} \max_{0 \leq \beta, \tau_h} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) := \frac{\beta \tau_g}{2} + L\left(\sqrt{\alpha^2 + a^2 \theta^2}, V\theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta}\right) \\ - \min_{\lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^q + \frac{\nu^2}{2} - \mathbb{E} \left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right); \lambda_0 \right) \right\} + \frac{\alpha \tau_h}{2} \right].$$

More specifically, a variety of conclusions can be derived based on the optimal solutions of the above optimization problem as we discuss in the next step. We note that while this expression may look complicated we prove that this optimization problem is in fact convex in the minimization parameters (θ, α, τ_g) and concave in the maximization parameters (β, τ_h) so that its optimal solutions can be easily derived via a simple low-dimensional gradient descent rather quickly and accurately. We also note that this proof is quite intricate and involved, so it is not possible to give an intuitive sketch of the arguments here. We refer to Section B for details. However, we briefly state a few mathematical challenges that is unique to simplifying (7.3). First, the AO (7.3) does not have a simple regularization whose scalarization reduces to a simple mean width calculation as in most simple CGMT uses. Instead the regularization has a complicated form $\frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q})$ which requires rather intricate and involved scalarization calculations. Second, this AO regularization term is not separable in $\boldsymbol{\theta}$ which significantly complicates the scalarization of the AO. Finally, we handle the case of more general covariance matrices where $\boldsymbol{\Sigma} \neq \mathbf{I}$.

Step IV: Completing the proof of the theorems.

Finally, we utilize the above scalar form to derive all of the different theorems and results. This is done by relating the quantities of interest in each theorem to the optimal solutions of (7.4). For instance, we show that $\lim_{n \rightarrow \infty} \text{SA}(\widehat{\boldsymbol{\theta}}^\varepsilon) = \Phi\left(\frac{V\theta_*}{\sqrt{\alpha_*^2 + a^2 \theta_*^2}}\right)$ and $\lim_{n \rightarrow \infty} \text{RA}(\widehat{\boldsymbol{\theta}}^\varepsilon) = \Phi\left(\frac{-\varepsilon_0 \gamma_{0*} + V\theta_*}{\sqrt{\alpha_*^2 + a^2 \theta_*^2}}\right)$ where α_* and θ_* are the optimal solutions over α and θ . These calculations/proofs are carried out in detail in Section B. Since each argument is different we do not provide a summary here and refer to the corresponding sections.

8. Discussion We conclude the paper by discussing some of the potential extensions and applications of our theory as well as comparison with more classical approaches to binary classification.

8.1. Generalization to random features models While our focus in this paper was on linear classifiers, these models are quite foundational and serve as

the basis for more complex models. For instance, one potential generalization of our results is to the class of random features models given by

$$\mathcal{F}_{\text{RF}} := \{f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{W}) = \text{sign}(\langle \boldsymbol{\theta}, \sigma(\mathbf{W}\mathbf{x}) \rangle) : \boldsymbol{\theta} \in \mathbb{R}^N\},$$

where $\mathbf{x} \in \mathbb{R}^d$ represents the feature vector, $\mathbf{W} \in \mathbb{R}^{N \times d}$ is a random matrix whose rows are chosen uniformly at random from the unit sphere in d -dimension, and $\sigma : \mathbb{R} \mapsto \mathbb{R}$ is a nonlinear function (for a vector \mathbf{v} , $\sigma(\mathbf{v}) = (\sigma(v_1), \dots, \sigma(v_m))$ is applied entry-wise). Random features model can also be described as a two-layer fully connected neural network with random first-layer weights fixed to \mathbf{W} and not optimized, while the second layer weights are represented by vector $\boldsymbol{\theta}$ and are optimized over to minimize the loss of interest. The random features model was introduced by [55] for scaling kernel methods to large datasets and there has been a large body of work drawing connections between random features models, kernel methods and fully trained neural networks [11, 10, 27, 37].

An intriguing phenomenon, pointed out by [45, 49] from the analysis of random features model in non-adversarial contexts, is that the random features model has the same asymptotic behavior as a simpler noisy linear features model whose second order statistics match the nonlinear random features model, namely a linear model with noisy features $\mathbf{u} \in \mathbb{R}^N$ given by $\mathbf{u} = \eta_0 + \eta_1 \mathbf{W}\mathbf{x} + \eta_2 \mathbf{z}$, where \mathbf{z} has i.i.d standard normal entries, independent of \mathbf{W} and \mathbf{x} . Also, the constants η_0, η_1, η_2 depend on the activation function $\sigma(\cdot)$ and are chosen so that the two models have the same first and second moments. A promising direction is to establish a similar connection for an adversarial setting and use our theory (relied on CGMT framework) to analyze the equivalent noisy linear model, from which we obtain an asymptotic characterization for adversarial training under the random features model. Very recently and after this paper was posted, [23] has pursued a similar approach to precisely characterize the role of overparametrization on robust generalization of random features in a regression setting.

8.2. Optimal ε_0 for the robust minimax estimator An interesting application of our theory is to derive the optimal value $\varepsilon_0^{\text{op}}$ (perceived perturbation level) in the robust minimax estimator (2.2), while fixing the adversary's (actual) perturbation level on test inputs to $\varepsilon_{0,\text{test}}$. (See Remark 4.4 on how our theory applies to this setting.) The optimality here is with respect to maximizing the robust accuracy. Somewhat surprisingly $\varepsilon_0^{\text{op}}$ is different than $\varepsilon_{0,\text{test}}$ in general and depends on δ and the choice of perturbation norm ℓ_p in a non-trivial way (There is no one-fit-all solution and this highlights the importance of having a precise theory to understand the effect of adversarial

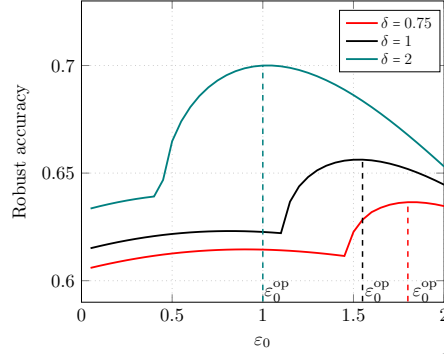


Fig 8: Robust accuracy curves versus ε_0 for different choices of δ , and the perturbation norm ℓ_1 ($p = 1$). The optimal choice of ε_0 for the robust minimax estimator decreases with δ .

training which is the primary goal of the current work). For example, in the particular case of $\varepsilon_{0,\text{test}} = 0$, the question reduces to finding the value of ε_0 which maximizes standard accuracy. As we already discussed, the answer very much depends on δ and p . For $p = 2$, we observe that (cf. Figure 3(a)) adversarial training helps with improving the standard accuracy. However for $p = \infty$, $\varepsilon_0^{\text{op}}$ should be large enough so that the problem becomes non-separable and also its value decreases as δ increases (cf. Figure 5(a)). As another example, we consider the case of $\varepsilon_{0,\text{test}} = 0.3$ with ℓ_∞ perturbations. In Figure 8 we plot the robust accuracy versus ε_0 , and the dashed vertical lines show the value of $\varepsilon_0^{\text{op}}$. As we see its value decreases by increasing δ , however, its exact value requires a precise analysis.

8.3. Comparison with Linear Discriminant Analysis (LDA) A classical approach to binary classification under the Gaussian-mixture model is the Linear Discriminant Analysis. In comparing the robustness property of LDA and the robust minimax estimator studied in this paper, we cannot say one estimator always outperforms the others. To further discuss this point, we consider the Gaussian-mixture model with identity covariance $\Sigma = \mathbf{I}$ and balanced classes. In this case, the LDA estimator reduces to $\hat{\mu}^{\text{LDA}} = \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i$ and the corresponding classification rule given by $\hat{y} = \text{sign}(\langle \mathbf{x}, \hat{\mu}^{\text{LDA}} \rangle)$. In the supplementary [29] (Section A), we compare the robust accuracy of LDA estimator with that of the robust minimax estimator $\hat{\theta}^\varepsilon$ for some choices of p . As we will discuss, the depending on p and the adversary's power ε_0 , one can outperform the other.

SUPPLEMENTARY MATERIAL

Supplement to: “Precise Statistical Analysis of Classification Accuracies for Adversarial Training”

(). Due to space constraints, proofs of theorems and some of the technical details are provided in the Supplementary Material [29].

References

- [1] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp. Living on the edge: A geometric theory of phase transitions in convex optimization. *arXiv preprint arXiv:1303.6672*, 2013.
- [2] D. Bartl, S. Drapeau, J. Obloj, and J. Wiesel. Robust uncertainty sensitivity analysis. *arXiv preprint arXiv:2006.12022*, 2020.
- [3] M. Bayati and A. Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *IEEE Transactions on Information Theory*, 57(2):764–785, 2011.
- [4] M. Bayati and A. Montanari. The lasso risk for gaussian matrices. *Information Theory, IEEE Transactions on*, 58(4):1997–2017, 2012.
- [5] M. Belkin, D. Hsu, S. Ma, and S. Mandal. Reconciling modern machine learning and the bias-variance trade-off. *arXiv preprint arXiv:1812.11118*, 2018.
- [6] M. Belkin, S. Ma, and S. Mandal. To understand deep learning we need to understand kernel learning. In *International Conference on Machine Learning*, pages 541–549, 2018.
- [7] Z. Bu, J. Klusowski, C. Rush, and W. Su. Algorithmic analysis and statistical estimation of slope via approximate message passing. In *Advances in Neural Information Processing Systems*, pages 9361–9371, 2019.
- [8] E. J. Candès and P. Sur. The phase transition for the existence of the maximum likelihood estimate in high-dimensional logistic regression. *The Annals of Statistics*, 48(1):27–42, 2020.
- [9] M. Celentano and A. Montanari. Fundamental barriers to high-dimensional regression with convex penalties. *arXiv preprint arXiv:1903.10603*, 2019.
- [10] A. Daniely. Sgd learns the conjugate kernel class of the network. In *Advances in Neural Information Processing Systems*, pages 2422–2430, 2017.
- [11] A. Daniely, R. Frostig, and Y. Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. In *Advances In Neural Information Processing Systems*, pages 2253–2261, 2016.
- [12] Z. Deng, A. Kammoun, and C. Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *arXiv preprint arXiv:1911.05822*, 2019.
- [13] Y. Deshpande and A. Montanari. Sparse pca via covariance thresholding. *Advances in Neural Information Processing Systems*, 27:334–342, 2014.
- [14] E. Dobriban, H. Hassani, D. Hong, and A. Robey. Provable tradeoffs in adversarially robust classification. *arXiv preprint arXiv:2006.05161*, 2020.
- [15] D. Donoho and A. Montanari. High dimensional robust m-estimation: Asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166(3-4):935–969, 2016.
- [16] D. L. Donoho, M. Gavish, and I. M. Johnstone. Optimal shrinkage of eigenvalues in the spiked covariance model. *Annals of statistics*, 46(4):1742, 2018.
- [17] D. L. Donoho, A. Javanmard, and A. Montanari. Information-theoretically optimal

- compressed sensing via spatial coupling and approximate message passing. *IEEE transactions on information theory*, 59(11):7434–7464, 2013.
- [18] D. L. Donoho, A. Maleki, and A. Montanari. The noise-sensitivity phase transition in compressed sensing. *Information Theory, IEEE Transactions on*, 57(10):6920–6941, 2011.
 - [19] N. El Karoui. On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields*, 170(1-2):95–175, 2018.
 - [20] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
 - [21] Y. Gordon. On milman’s inequality and random subspaces which escape through a mesh in r^n . In *Geometric aspects of functional analysis*, pages 84–106. Springer, 1988.
 - [22] M. Grant, S. Boyd, and Y. Ye. *Cvx: Matlab software for disciplined convex programming*, 2008.
 - [23] H. Hassani and A. Javanmard. The curse of overparametrization in adversarial training: Precise analysis of robust generalization for random features regression. *arXiv preprint arXiv:2201.05149*, 2022.
 - [24] T. Hastie, A. Montanari, S. Rosset, and R. J. Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
 - [25] H. Hu and Y. M. Lu. Asymptotics and optimal designs of slope for sparse linear regression. *arXiv preprint arXiv:1903.11582*, 2019.
 - [26] H. Huang. Asymptotic behavior of support vector machine for spiked population model. *The Journal of Machine Learning Research*, 18(1):1472–1492, 2017.
 - [27] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
 - [28] A. Javanmard, A. Montanari, and F. Ricci-Tersenghi. Phase transitions in semidefinite relaxations. *Proceedings of the National Academy of Sciences*, 113(16):E2218–E2223, 2016.
 - [29] A. Javanmard and M. Soltanolkotabi. Supplementary material to “precise statistical analysis of classification accuracies for adversarial training”. 2020.
 - [30] A. Javanmard, M. Soltanolkotabi, and H. Hassani. Precise tradeoffs in adversarial training for linear regression. volume 125 of *Proceedings of Machine Learning Research, Conference of Learning Theory (COLT)*, pages 2034–2078. PMLR, 09–12 Jul 2020.
 - [31] Z. Ji and M. Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018.
 - [32] I. M. Johnstone. On the distribution of the largest eigenvalue in principal components analysis. *Annals of statistics*, pages 295–327, 2001.
 - [33] A. Kammoun and M.-S. Alouini. On the precise error analysis of support vector machines. *Submitted to IEEE Transactions on information theory*, 2019.
 - [34] N. E. Karoui. Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results. *arXiv preprint arXiv:1311.2445*, 2013.
 - [35] A. Kurakin, I. Goodfellow, and S. Bengio. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*, 2016.
 - [36] L. Lai and E. Bayraktar. On the adversarial robustness of robust estimators. *IEEE Transactions on Information Theory*, 66(8):5097–5109, 2020.
 - [37] Y. Li and Y. Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. *NeurIPS*, 2018.

- [38] T. Liang and P. Sur. A precise high-dimensional asymptotic theory for boosting and min-l1-norm interpolated classifiers. *arXiv preprint arXiv:2002.01586*, 2020.
- [39] F. Liese and K.-J. Miescke. Statistical decision theory: Estimation, testing, and selection. In *Springer Science & Business Media*, 2008.
- [40] P. Lolas. Regularization in high-dimensional regression and classification via random matrix theory. *arXiv preprint arXiv:2003.13723*, 2020.
- [41] K. Lyu and J. Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- [42] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- [43] X. Mai, Z. Liao, and R. Couillet. 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)*, pages 3357–3361. IEEE, 2019.
- [44] M. Mehrabi, A. Javanmard, R. A. Rossi, A. Rao, and T. Mai. Fundamental trade-offs in distributionally adversarial training. volume 139 of *Proceedings of the 38th International Conference on Machine Learning*, pages 7544–7554. PMLR, 2021.
- [45] S. Mei and A. Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- [46] F. Mignacco, F. Krzakala, Y. M. Lu, and L. Zdeborová. The role of regularization in classification of high-dimensional noisy gaussian mixture. *arXiv preprint arXiv:2002.11544*, 2020.
- [47] Y. Min, L. Chen, and A. Karbasi. The curious case of adversarially robust models: More data can help, double descend, or hurt generalization. *arXiv preprint arXiv:2002.11080*, 2020.
- [48] L. Miolane and A. Montanari. The distribution of the lasso: Uniform control over sparse balls and adaptive parameter tuning. *arXiv preprint arXiv:1811.01212*, 2018.
- [49] A. Montanari, F. Ruan, Y. Sohn, and J. Yan. The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime. *arXiv preprint arXiv:1911.01544*, 2019.
- [50] S. Oymak, B. Recht, and M. Soltanolkotabi. Sharp time–data tradeoffs for linear inverse problems. *IEEE Transactions on Information Theory*, 64(6):4129–4158, 2017.
- [51] S. Oymak, C. Thrampoulidis, and B. Hassibi. The squared-error of generalized lasso: A precise analysis. *arXiv preprint arXiv:1311.0830*, 2013.
- [52] M. S. Pydi and V. Jog. Adversarial risk via optimal transport and optimal couplings. In *International Conference on Machine Learning*, pages 7814–7823. PMLR, 2020.
- [53] A. Raghunathan, J. Steinhardt, and P. Liang. Certified defenses against adversarial examples. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- [54] A. Raghunathan, S. M. Xie, F. Yang, J. C. Duchi, and P. Liang. Adversarial training can hurt generalization. *arXiv preprint arXiv:1906.06032*, 2019.
- [55] A. Rahimi and B. Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20:1177–1184, 2007.
- [56] F. Salehi, E. Abbasi, and B. Hassibi. The impact of regularization on high-dimensional logistic regression. *arXiv preprint arXiv:1906.03761*, 2019.
- [57] L. Schmidt, S. Santurkar, D. Tsipras, K. Talwar, and A. Madry. Adversarially robust generalization requires more data. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*, pages 5019–5031, 2018.

- [58] M. Staib and S. Jegelka. Distributionally robust deep learning as a generalization of adversarial training. In *NIPS workshop on Machine Learning and Computer Security*, 2017.
- [59] M. Stojnic. Various thresholds for ℓ_1 -optimization in compressed sensing. *arXiv preprint arXiv:0907.3666*, 2009.
- [60] M. Stojnic. A framework to characterize performance of lasso algorithms. *arXiv preprint arXiv:1303.7291*, 2013.
- [61] P. Sur and E. J. Candès. A modern maximum-likelihood theory for high-dimensional logistic regression. *Proceedings of the National Academy of Sciences*, 116(29):14516–14525, 2019.
- [62] H. Taheri, R. Pedarsani, and C. Thrampoulidis. Sharp asymptotics and optimal performance for inference in binary models. *arXiv preprint arXiv:2002.07284*, 2020.
- [63] C. Thrampoulidis, E. Abbasi, and B. Hassibi. Precise error analysis of regularized m -estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628, 2018.
- [64] C. Thrampoulidis, S. Oymak, and B. Hassibi. Regularized linear regression: A precise analysis of the estimation error. In *Proceedings of The 28th Conference on Learning Theory*, pages 1683–1709, 2015.
- [65] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- [66] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.
- [67] S. Wang, H. Weng, and A. Maleki. Does slope outperform bridge regression? *arXiv preprint arXiv:1909.09345*, 2019.
- [68] E. Wong and J. Z. Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pages 5283–5292, 2018.
- [69] H. Zhang, Y. Yu, J. Jiao, E. P. Xing, L. E. Ghaoui, and M. I. Jordan. Theoretically principled trade-off between robustness and accuracy. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pages 7472–7482, 2019.

SUPPLEMENTARY MATERIAL TO “PRECISE STATISTICAL ANALYSIS OF CLASSIFICATION ACCURACIES FOR ADVERSARIAL TRAINING”

BY ADEL JAVANMARD[§] AND MAHDI SOLTANOLKOTABI[§]

University of Southern California[§]

APPENDIX A: COMPARISON WITH LINEAR DISCRIMINANT ANALYSIS (LDA)

A classical approach to binary classification under the Gaussian-mixture model is the Linear Discriminant Analysis. In comparing the robustness property of LDA and the robust minimax estimator studied in this paper, we cannot say one estimator always outperforms the others. To further discuss this point, we consider the Gaussian-mixture model with identity covariance $\Sigma = \mathbf{I}$ and balanced classes. In this case, the LDA estimator reduces to $\hat{\boldsymbol{\mu}}^{\text{LDA}} = \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i$ and the corresponding classification rule given by $\hat{y} = \text{sign}(\langle \mathbf{x}, \hat{\boldsymbol{\mu}}^{\text{LDA}} \rangle)$. Under the Gaussian-mixture model we have $\mathbf{x} = y\boldsymbol{\mu} + \mathbf{z}$ with $\mathbf{z} \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$. Therefore,

$$\hat{\boldsymbol{\mu}}^{\text{LDA}} = \frac{1}{n} \sum_{i=1}^n y_i (y_i \boldsymbol{\mu} + \mathbf{z}_i) = \boldsymbol{\mu} + \frac{1}{n} \sum_{i=1}^n y_i \mathbf{z}_i = \boldsymbol{\mu} + \tilde{\mathbf{z}}, \quad \tilde{\mathbf{z}} \sim \mathbf{N}(\mathbf{0}, \frac{1}{n} \mathbf{I})$$

For simplicity we assume that the class averages $\boldsymbol{\mu}$ is generated as $\boldsymbol{\mu} \sim (\mathbf{0}, \frac{1}{d} \mathbf{I})$, similar to the setting considered in the numerical experiments. In asymptotic regime of $n \rightarrow \infty$ and $n/d \rightarrow \delta$, we have that in probability:

$$\begin{aligned} \lim_{n \rightarrow \infty} \langle \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}^{\text{LDA}} \rangle &= \lim_{n \rightarrow \infty} \|\boldsymbol{\mu}\|_{\ell_2}^2 = 1, \\ \lim_{n \rightarrow \infty} d^{1/2-1/q} \|\hat{\boldsymbol{\mu}}^{\text{LDA}}\|_{\ell_q} &= \lim_{n \rightarrow \infty} d^{1/2-1/q} \|\boldsymbol{\mu} + \tilde{\mathbf{z}}\|_{\ell_q} \\ &= \lim_{n \rightarrow \infty} d^{1/2-1/q} \left(\frac{1}{d} + \frac{1}{n} \right)^{1/2} d^{1/q} C_q = \left(1 + \frac{1}{\delta} \right)^{1/2} C_q, \end{aligned}$$

where in the first equation we used the fact that $\langle \boldsymbol{\mu}, \tilde{\mathbf{z}} \rangle \sim \mathbf{N}(0, \frac{1}{n} \|\boldsymbol{\mu}\|_{\ell_2}^2)$ has vanishing variance as $n \rightarrow \infty$. In the second inequality, C_q is the q -th moment of standard normal distribution. Recall that $\varepsilon = \varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}$ with $1/p + 1/q = 1$, and also $\|\boldsymbol{\mu}\|_{\ell_p} \rightarrow d^{1/p-1/2} C_p = d^{1/2-1/q} C_p$. Using these identities along with

the characterization of standard and robust accuracies given by Lemma 2.1 of the paper, we arrive at

$$(A.1) \quad \begin{aligned} \lim_{n \rightarrow \infty} SA(\widehat{\boldsymbol{\mu}}^{\text{LDA}}) &= \Phi\left(\sqrt{\frac{\delta}{1+\delta}}\right), \\ \lim_{n \rightarrow \infty} RA(\widehat{\boldsymbol{\mu}}^{\text{LDA}}) &= \Phi\left(\sqrt{\frac{\delta}{1+\delta}} - \varepsilon_0 C_q C_p\right). \end{aligned}$$

We next compare the robust accuracy of LDA estimator with that of the robust minimax estimator $\widehat{\boldsymbol{\theta}}^\varepsilon$ for some choices of p . As we will discuss, the depending on p and the adversary's power ε_0 , one can outperform the other.

- ($p = q = 2$). Figure 9(a) compares $RA(\widehat{\boldsymbol{\mu}}^{\text{LDA}})$ with $RA(\widehat{\boldsymbol{\theta}}^\varepsilon)$ versus ε_0 for several values of δ . Here, the solid lines correspond to the robust minimax estimator and the dashed lines correspond to the LDA estimator. Figure 9(b) compares $RA(\widehat{\boldsymbol{\mu}}^{\text{LDA}})$ with $RA(\widehat{\boldsymbol{\theta}}^\varepsilon)$ versus $1/\delta$ for various choices of ε_0 . As we see for the case of $p = 2$, the LDA has better robust accuracy and it is mostly very close to that of the robust estimator.
- ($p = \infty, q = 1$). Similar to the setting of experiments in Section 5.2, here we consider the scaling $\varepsilon = \varepsilon_0/\sqrt{d}$. Figure 10 (a) compares the robust accuracies versus ε_0 for several values of δ . As we see for any δ , there exists $\varepsilon_0^*(\delta)$ above which the robust minimax outperforms the LDA. Figure 10(b) compares the robust accuracies versus $1/\delta$ for several values of ε_0 . Rewording the above observation, for any ε_0 there exists $\delta^*(\varepsilon_0)$ below which the robust minimax outperforms the LDA estimator.
- ($p = 1, q = \infty$). Similar to the setting of experiments in Section 5.3, we have $\varepsilon = \varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p} = \sqrt{\frac{2}{\pi}} \frac{\varepsilon_0}{\sqrt{d}}$. Invoking equations (A.1), we have $\lim_{n \rightarrow \infty} RA(\widehat{\boldsymbol{\mu}}^{\text{LDA}}) = 0$ because $C_q = \sqrt{2 \log d} \rightarrow \infty$. However, as we see in Figure 7, the robust minimax estimator $\widehat{\boldsymbol{\theta}}^\varepsilon$ achieves non-trivial positive robust accuracies and hence outperforms LDA.

APPENDIX B: PROOFS FOR ANISOTROPIC GAUSSIAN MODEL (SECTION 6)

B.1. Proof of Theorem 6.1 As discussed the (ε, q) -separability condition can alternatively be written as (3.2), which we repeat here:

$$(B.1) \quad \exists \boldsymbol{\theta}, \quad \|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon} : \quad \forall i \in [n], \quad y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle > 1.$$

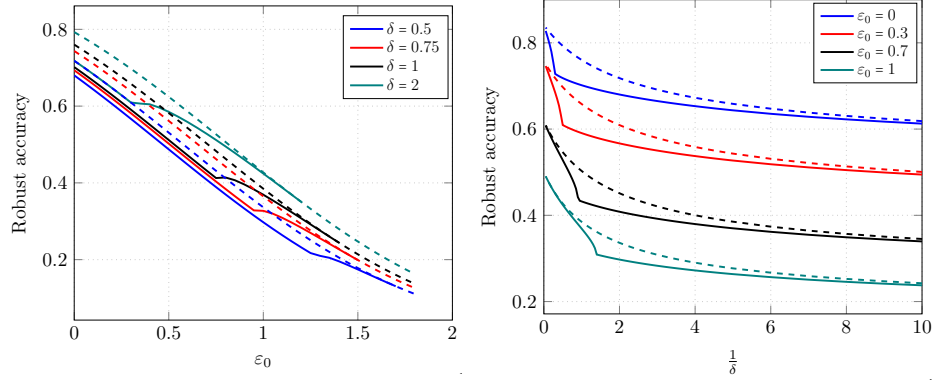


Fig 9: Robust accuracies for the LDA estimator and the robust minimax estimator versus the adversary's power with ℓ_2 ($p = 2$) perturbations for different values of δ . Solid curves correspond to the robust minimax estimator and the dashed curves correspond to the LDA estimator.

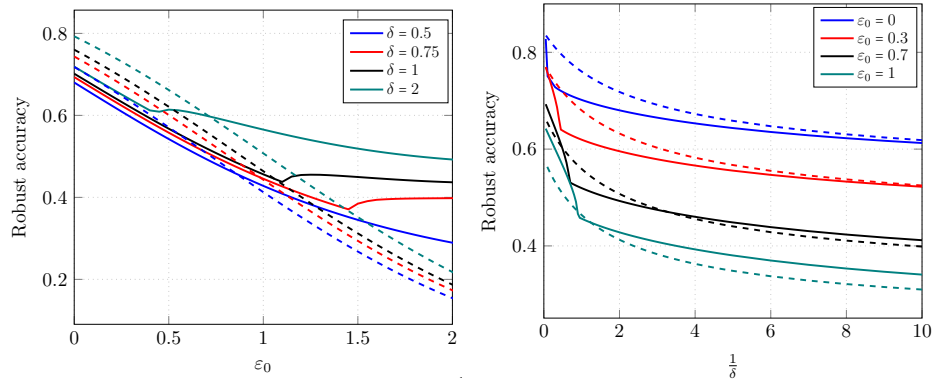


Fig 10: Robust accuracies for LDA estimator and the robust minimax estimator versus the adversary's power with ℓ_∞ ($p = \infty$) perturbations for different values of δ . Solid curves correspond to the robust minimax estimator and the dashed curves correspond to the LDA estimator.

To find the separability threshold we consider the following feasibility problem

$$(B.2) \quad \min_{\boldsymbol{\theta} \in \mathbb{R}^d} 0 \quad \text{subject to } y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle > 1, \quad \|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon}.$$

Clearly this is a convex optimization problem since $q \geq 1$. Writing the partial Lagrangian for the above problem with u_i/n as dual coefficients, this is equivalent to

$$(B.3) \quad \min_{\boldsymbol{\theta}} \max_{u_i \geq 0} \frac{1}{n} \sum_{i=1}^n u_i (1 - y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle) \quad \text{subject to } \|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon}.$$

Under our Gaussian Mixture data model, we can substitute for $\mathbf{X} = \mathbf{y}\boldsymbol{\mu}^T + \mathbf{Z}\boldsymbol{\Sigma}^{1/2}$, which results in

$$(B.4) \quad \min_{\boldsymbol{\theta}} \max_{u_i \geq 0} \frac{1}{n} \mathbf{u}^T \left(\mathbf{1} (1 - \boldsymbol{\mu}^T \boldsymbol{\theta}) - \mathbf{D}_y \mathbf{Z} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right) \quad \text{subject to } \|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon}.$$

The above dual problem has a finite optimal value if and only if the data is (ε, q) -separable. So we aim at finding the largest δ such that the above problem has still a finite optimal value. (Recall that $\frac{n}{d} \rightarrow \delta$.)

Reduction to an auxiliary optimization problem via CGMT. Note that $y_i = \pm 1$ are independent of \mathbf{Z} . In addition, the objective function in (B.4) is affine in the standard Gaussian matrix \mathbf{Z} and the rest of the terms form a convex-concave function in $\boldsymbol{\theta}, \mathbf{u}$. Due to this particular form we are able to apply a powerful extension of a classical Gaussian process inequality due to Gordon [21] known as Convex Gaussian Minimax Theorem (CGMT) [64], and is discussed in the proof sketch in Section 7. The CGMT framework provides a principled machinery to characterize the asymptotic behavior of certain minimax optimization problems that are affine in a Gaussian matrix \mathbf{X} .

As discussed in the CGMT framework in Section 7, we require minimization/maximization to be over compact sets. The vector $\boldsymbol{\theta}$ already lies in the ℓ_q ball of radius $1/\varepsilon$ by constraint. In addition, since $u_i \geq 0$, and we are focused on the regime that (B.4) has finite optimal value, the optimal values of u_i should all be finite as well.

We are now ready to applying the CGMT framework. The corresponding Auxiliary Optimization (AO) reads as

$$(B.5) \quad \min_{\boldsymbol{\theta}} \max_{\mathbf{u} \geq 0} \frac{1}{n} \left\{ (\mathbf{u}^T \mathbf{1}) (1 - \boldsymbol{\mu}^T \boldsymbol{\theta}) + \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g}^T \mathbf{u} + \|\mathbf{u}\|_{\ell_2} \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\},$$

$$\text{subject to } \|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon},$$

where $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_n)$ and $\mathbf{h} \sim \mathcal{N}(0, \mathbf{I}_d)$. Fixing $\beta := \frac{\|\mathbf{u}\|_{\ell_2}}{\sqrt{n}}$ and optimizing over \mathbf{u} on the non-negative orthant we get

$$\min_{\boldsymbol{\theta}} \max_{\beta \geq 0} \quad \frac{\beta}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + \frac{\beta}{\sqrt{n}} \left\| \left((1 - \boldsymbol{\mu}^T \boldsymbol{\theta}) \mathbf{1} + \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} \right)_+ \right\|_{\ell_2},$$

(B.6) subject to $\|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon}.$

For data to be separable the above dual optimization should take finite optimal value and therefore the coefficient of β should be non-positive. As such the problem is separable if and only if the optimal value of the following problem is non-positive:

$$\min_{\boldsymbol{\theta}} \frac{1}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + \frac{1}{\sqrt{n}} \left\| \left((1 - \boldsymbol{\mu}^T \boldsymbol{\theta}) \mathbf{1} + \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} \right)_+ \right\|_{\ell_2} \leq 0,$$

(B.7) subject to $\|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon}.$

Consider the decomposition $\boldsymbol{\theta} = \boldsymbol{\theta}_\perp + \theta \tilde{\boldsymbol{\mu}}$ with $\boldsymbol{\theta}_\perp = \mathbf{P}_\mu^\perp \boldsymbol{\theta}$. Note that

$$\begin{aligned} \frac{1}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} &= \frac{1}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_\perp + \frac{1}{\sqrt{n}} \theta \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} \\ &= \frac{1}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_\perp + \frac{1}{\sqrt{n}} \theta \mathbf{h}^T \tilde{\boldsymbol{\mu}} \end{aligned}$$

Since $\mathbf{h}^T \tilde{\boldsymbol{\mu}} \sim \mathcal{N}(0, 1)$ and θ is bounded the contribution of the second term is negligible in the large sample limit $n \rightarrow \infty$. This along with the symmetry of the distribution of \mathbf{h} bring us to

$$\min_{\alpha \geq 0, \theta, \boldsymbol{\theta}_\perp} -\frac{1}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_\perp + \frac{1}{\sqrt{n}} \left\| \left((1 - \|\boldsymbol{\mu}\|_{\ell_2} \theta) \mathbf{1} + \sqrt{\alpha^2 + a^2 \theta^2} \mathbf{g} \right)_+ \right\|_{\ell_2}$$

(B.8) subject to $\|\boldsymbol{\theta}\|_{\ell_q} \leq \frac{1}{\varepsilon}, \quad \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_\perp \right\|_{\ell_2} = \alpha, \quad \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} = \theta$

Scalarization of the auxiliary optimization problem. To continue recall the definition of set \mathcal{S} given by

$$\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu}) := \left\{ \mathbf{z} \in \mathbb{R}^d : \mathbf{z}^T \tilde{\boldsymbol{\mu}} = 0, \|\mathbf{z}\|_{\ell_2} = \alpha, \left\| \boldsymbol{\Sigma}^{-1/2} \mathbf{z} + \theta \tilde{\boldsymbol{\mu}} \right\|_{\ell_q} \leq \frac{1}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}} \right\}.$$

Recall that $\varepsilon = \varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}$ and so the optimization problem (B.8) above can be rewritten in the form

(B.9)

$$\min_{\alpha \geq 0, \theta} \min_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})} -\frac{1}{\sqrt{n}} \mathbf{h}^T \mathbf{z} + \frac{1}{\sqrt{n}} \left\| \left((1 - \|\boldsymbol{\mu}\|_{\ell_2} \theta) \mathbf{1} + \sqrt{\alpha^2 + a^2 \theta^2} \mathbf{g} \right)_+ \right\|_{\ell_2}$$

Recall the spherical width of a set $\mathcal{S} \subset \mathbb{R}^d$ defined as

$$\omega_s(\mathcal{S}) = \mathbb{E} \left[\sup_{\mathbf{z} \in \mathcal{S}} \mathbf{z}^T \mathbf{u} \right],$$

where $\mathbf{u} \in \mathbb{S}^{d-1}$ is a vector chosen at random from the unit sphere. Using this definition and the fact that $\min z = -\max -z$ we have

$$\min_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0)} -\frac{1}{\sqrt{n}} \mathbf{h}^T \mathbf{z} = -\frac{1}{\sqrt{n/d}} \sup_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0)} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} \rightarrow -\frac{1}{\sqrt{\delta}} \omega(\alpha, \theta, \varepsilon_0),$$

in probability, where in the last line we use the fact, for $\mathcal{S} \in \mathbb{S}^{d-1}$, the function $f(\mathbf{u}) = \sup_{\mathbf{z} \in \mathcal{S}} \mathbf{z}^T \mathbf{u}$ is Lipschitz. Therefore, using the concentration of Lipschitz functions of Gaussian random vectors (see e.g. [66, Theorem 5.2.2]), $f(\mathbf{u})$ concentrates around its mean $\mathbb{E} f(\mathbf{u}) = \omega_s(\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu}))$. More precisely,

$$\mathbb{P} \left\{ \left| \sup_{\mathbf{z} \in \mathcal{S}} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} - \omega_s(\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})) \right| \right\} \leq 2e^{-cdt^2},$$

for an absolute constant $c > 0$ and for every $t \geq 0$. Therefore, by invoking the assumption on the convergence of spherical width, cf. Assumption 3, we arrive at

$$\lim_{d \rightarrow \infty} \mathbb{P} \left\{ \left| \sup_{\mathbf{z} \in \mathcal{S}} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} - \omega(\alpha, \theta, \varepsilon_0) \right| \geq \eta \right\} = 0, \quad \forall \eta > 0.$$

Therefore, $\sup_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0)} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} \rightarrow \omega(\alpha, \theta, \varepsilon_0)$, in probability.

Furthermore, $\|\boldsymbol{\mu}\|_{\ell_2} \rightarrow V$ by Assumption 2 and since $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_n)$ by applying the Weak Law of Large Numbers we have

$$\frac{1}{\sqrt{n}} \left\| \left((1 - \|\boldsymbol{\mu}\|_{\ell_2} \theta) \mathbf{1} + \sqrt{\alpha^2 + a^2 \theta^2} \mathbf{g} \right)_+ \right\|_{\ell_2} \rightarrow \sqrt{\mathbb{E} \left[\left((1 - V\theta + \sqrt{\alpha^2 + a^2 \theta^2} g \right)_+^2 \right]}$$

Thus the objective function in the optimization problem (B.9) converges pointwise to

$$(B.10) \quad \min_{\alpha \geq 0, \theta} -\frac{1}{\sqrt{\delta}} \omega(\alpha, \theta, \varepsilon_0) + \sqrt{\mathbb{E} \left[\left((1 - V\theta + \sqrt{\alpha^2 + a^2 \theta^2} g \right)_+^2 \right]}$$

Also the problem (B.9) is convex as a function of $(\alpha, \theta, \mathbf{z})$ and since partial maximization preserves convexity, the objective of (B.9) (after minimization over \mathbf{z}) is a convex function of (α, θ) . We can thus apply the convexity lemma [63, Lemma B.2] to conclude that the minimum value of (B.9) over

$\alpha \geq 0, \theta$ also converges to that of (B.10). Therefore, we conclude that data is (ε, q) -separable if and only if the optimal value in (B.10) is finite. Rearranging the terms gives us that (B.10) has a finite optimal value if and only if

$$(B.11) \quad \delta < \delta_*, \quad \text{with } \delta_* := \min_{\alpha \geq 0, \theta} \frac{\omega(\alpha, \theta, \varepsilon_0)^2}{\mathbb{E} \left[\left(1 - V\theta + \sqrt{\alpha^2 + \theta^2} g \right)_+^2 \right]}.$$

This completes the proof of Theorem 6.1.

B.2. Proof of Theorem 6.3 We prove Theorem 6.3 using the Convex Gaussian Minimax Theorem (CGMT) as outlined in Section 7. The max-margin problem (F.2) can be equivalently written as

$$(B.12) \quad (\tilde{\theta}^\varepsilon, \hat{\gamma}) = \arg \min_{\theta, \gamma \geq 0} \|\theta\|_{\ell_2}^2 \\ \text{subject to } y_i \langle \mathbf{x}_i, \theta \rangle - \varepsilon \gamma \geq 1, \quad \gamma \geq \|\theta\|_{\ell_q}$$

Now note that writing the Lagrangian for the max-margin problem with u_i/n and 2λ as dual coefficients, this is equivalent to

$$(B.13) \quad \min_{\theta, \gamma \geq 0} \max_{u_i, \lambda \geq 0} \|\theta\|_{\ell_2}^2 + \frac{1}{n} \sum_{i=1}^n u_i (1 + \varepsilon \gamma - y_i \langle \mathbf{x}_i, \theta \rangle) + 2\lambda (\|\theta\|_{\ell_q} - \gamma).$$

We next substitute for $\mathbf{X} = \mathbf{y}\boldsymbol{\mu}^T + \mathbf{Z}\boldsymbol{\Sigma}^{1/2}$ based on the Gaussian mixtures model to arrive at

$$(B.14) \quad \min_{\theta, \gamma \geq 0} \max_{u_i \geq 0, \lambda \geq 0} \|\theta\|_{\ell_2}^2 + \frac{1}{n} \left(\mathbf{u}^T \mathbf{1} + \varepsilon \gamma \mathbf{u}^T \mathbf{1} - \mathbf{u}^T \mathbf{D}_y \mathbf{Z} \boldsymbol{\Sigma}^{1/2} \theta - \mathbf{u}^T \mathbf{1} \boldsymbol{\mu}^T \theta \right) + 2\lambda (\|\theta\|_{\ell_q} - \gamma).$$

The advantage of the Lagrangian form in (B.14) is that it is a minimax problem and the objective is an affine function of the standard Gaussian matrix \mathbf{Z} . Therefore, we can deploy the Convex Gaussian Minimax Theorem (CGMT) [64], described in Section 7, to characterize asymptotic values of certain functions of this optimization solution, in a high probability sense.

To recall, the CGMT framework shows that a problem of the form

$$(B.15) \quad \min_{\theta \in \mathcal{S}_\theta} \max_{\mathbf{u} \in \mathcal{S}_u} \mathbf{u}^T \mathbf{Z} \theta + \psi(\theta, \mathbf{u})$$

with \mathbf{Z} a matrix with $\mathcal{N}(0, 1)$ entries can be replaced asymptotically with

$$(B.16) \quad \min_{\theta \in \mathcal{S}_\theta} \max_{\mathbf{u} \in \mathcal{S}_u} \|\theta\|_{\ell_2} \mathbf{g}^T \mathbf{u} + \|\mathbf{u}\|_{\ell_2} \mathbf{h}^T \theta + \psi(\theta, \mathbf{u})$$

where \mathbf{g} and \mathbf{h} are independent Gaussian vectors with i.i.d. $\mathcal{N}(0, 1)$ entries and $\psi(\boldsymbol{\theta}, \mathbf{u})$ is convex in $\boldsymbol{\theta}$ and concave in \mathbf{u} . Specifically, the optimal value and corresponding solution of (B.15) converge in probability to the optimal value and the corresponding solution of (B.16). In the above \mathcal{S}_θ and \mathcal{S}_u are compact sets. We refer to [64, Theorem 3] for precise statements. As explained in the proof sketch in 7, we follow [64] in referring to problems of the form (B.15) and (B.16) as the Primal Problem (PO) and the Auxiliary Problem (AO).

Note that in order to apply CGMT, we need the minimization/maximization to be over compact sets. This technical issue can be avoided by introducing “artificial” boundedness constraints on the optimization variables that they do not change the optimal solution. Concretely, we can add constraints of the form $\mathcal{S}_\theta = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_{\ell_q} \leq K_\theta\}$ and $\mathcal{S}_u = \{\mathbf{u} : 0 \leq u_i, \frac{1}{n}\mathbf{1}^T \mathbf{u} \leq K_u\}$ for sufficiently large constants K_θ, K_u without changing the optimal solution of (B.14) in a precise asymptotic sense. We refer to Appendix E.3.1 for precise statements and proofs. This allows us to replace (B.14) with

$$(B.17) \quad \min_{\boldsymbol{\theta} \in \mathcal{S}_\theta, \gamma \geq 0} \max_{\mathbf{u} \in \mathcal{S}_u, \lambda \geq 0} \|\boldsymbol{\theta}\|_{\ell_2}^2 + \frac{1}{n} \left(\mathbf{u}^T \mathbf{1} + \varepsilon \gamma \mathbf{u}^T \mathbf{1} - \mathbf{u}^T \mathbf{D}_y \mathbf{Z} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} - \mathbf{u}^T \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} \right) + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right).$$

Reduction to an auxiliary optimization problem via CGMT. With these compact constraints in place we can now apply the CGMT result to obtain the auxiliary optimization (AO) problem.

We proceed by defining the projection matrices

$$\mathbf{P}_\mu^\perp := \mathbf{I} - \tilde{\boldsymbol{\mu}} \tilde{\boldsymbol{\mu}}^T, \quad \mathbf{P}_\mu := \tilde{\boldsymbol{\mu}} \tilde{\boldsymbol{\mu}}^T$$

and rewrite $\mathbf{Z} \boldsymbol{\Sigma}^{1/2} = \mathbf{Z} (\mathbf{P}_\mu + \mathbf{P}_\mu^\perp) \boldsymbol{\Sigma}^{1/2}$. Since $\mathbf{Z} \mathbf{P}_\mu$ and $\mathbf{Z} \mathbf{P}_\mu^\perp$ are independent from each other the latter has the same distribution as

$$\mathbf{Z} \boldsymbol{\Sigma}^{1/2} \sim \mathbf{z} \left(\boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} \right)^T + \mathbf{Z} \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2}$$

where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_n)$ and is independent from the matrix \mathbf{Z} . This brings us to the following representation

$$\begin{aligned} \min_{\boldsymbol{\theta}, \gamma \geq 0} \max_{\mathbf{u} \geq 0, \lambda \geq 0} & \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) \\ & + \frac{1}{n} \left\{ \mathbf{u}^T \mathbf{1} + \varepsilon \gamma \mathbf{u}^T \mathbf{1} - \mathbf{u}^T \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} - (\mathbf{u}^T \mathbf{D}_y \mathbf{z}) (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) - \mathbf{u}^T \mathbf{D}_y \mathbf{Z} \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\} \end{aligned}$$

Since $y_i = \pm 1$ are independent of \mathbf{Z} , by applying CGMT framework, the AO reads as

$$\begin{aligned} \min_{\boldsymbol{\theta}, \gamma \geq 0} \max_{\mathbf{u} \geq 0, \lambda \geq 0} \quad & \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) + \frac{1}{n} \left\{ \mathbf{u}^T \mathbf{1} + \varepsilon \gamma \mathbf{u}^T \mathbf{1} - \mathbf{u}^T \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} \right. \\ & + (\mathbf{u}^T \mathbf{z})(\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) + \left\| \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g}^T \mathbf{u} \\ & \left. + \|\mathbf{u}\|_{\ell_2} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\} \end{aligned}$$

Fixing $\beta := \frac{\|\mathbf{u}\|_{\ell_2}}{\sqrt{n}}$ and optimizing over \mathbf{u} on the non-negative orthant we get

$$\begin{aligned} \min_{\boldsymbol{\theta}, \gamma \geq 0} \max_{\beta \geq 0, \lambda \geq 0} \quad & \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \\ \text{(B.18)} \quad & + \frac{\beta}{\sqrt{n}} \left\| \left((1 + \varepsilon \gamma - \boldsymbol{\mu}^T \boldsymbol{\theta}) \mathbf{1} + (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) \mathbf{z} + \left\| \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} \right)_+ \right\|_{\ell_2} \end{aligned}$$

Since $\mathbf{z}, \mathbf{g} \sim \mathbf{N}(0, \mathbf{I}_n)$ are independent, by applying the Weak Law of Large Numbers we have

$$\begin{aligned} & \frac{1}{\sqrt{n}} \left\| \left((1 + \varepsilon \gamma - \boldsymbol{\mu}^T \boldsymbol{\theta}) \mathbf{1} + (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) \mathbf{z} + \left\| \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} \right)_+ \right\|_{\ell_2} \\ & \rightarrow \left(\mathbb{E} \left[\left((1 + \varepsilon \gamma - \boldsymbol{\mu}^T \boldsymbol{\theta}) + \sqrt{(\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta})^2 + \left\| \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2}^2} \mathbf{g} \right)_+^2 \right] \right)^{\frac{1}{2}} \\ & = \left(\mathbb{E} \left[\left((1 + \varepsilon \gamma - \boldsymbol{\mu}^T \boldsymbol{\theta}) + \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} \right)_+^2 \right] \right)^{\frac{1}{2}} \end{aligned}$$

Thus we arrive at

$$\begin{aligned} \min_{\boldsymbol{\theta}, \gamma \geq 0} \max_{\beta \geq 0, \lambda \geq 0} \quad & \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \\ \text{(B.19)} \quad & + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon \gamma - \boldsymbol{\mu}^T \boldsymbol{\theta}) + \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} \right)_+^2 \right]}. \end{aligned}$$

We note that for $a \geq 0$,

$$\mathbb{E}[ag + b]_+^2 = \frac{a^2 + b^2}{2} \left(1 + \operatorname{erf} \left(\frac{b}{\sqrt{2}a} \right) \right) + \frac{ab}{\sqrt{2\pi}} e^{-\frac{b^2}{2a^2}}.$$

and its derivative with respect to a is given by $2a(1 + \operatorname{erf}(\frac{b}{\sqrt{2}a})) > 0$ which implies that the function is increasing in $a > 0$. Therefore the optimization

(B.19) can be equivalently written as

$$\begin{aligned} & \min_{\boldsymbol{\theta}, \gamma, \alpha \geq 0} \max_{\beta \geq 0, \lambda \geq 0} \quad \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon \gamma - \boldsymbol{\mu}^T \boldsymbol{\theta}) + \alpha g \right)_+^2 \right]} \\ & \text{subject to} \quad \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2} \leq \alpha. \end{aligned} \quad (\text{B.20})$$

Note that the above is trivially jointly convex in $(\boldsymbol{\theta}, \gamma, \alpha)$ and jointly concave in (β, λ) . We fix the parallel component of $\boldsymbol{\theta}$ on $\boldsymbol{\mu}$ to θ , namely $\boldsymbol{\theta} = \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta}$. We next optimize over $\boldsymbol{\theta}$ while fixing θ .

$$\begin{aligned} & \min_{\boldsymbol{\theta}, \gamma \geq 0, \alpha \geq 0} \max_{\beta \geq 0, \lambda \geq 0} \quad \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon \gamma - \theta \|\boldsymbol{\mu}\|_{\ell_2}) + \alpha g \right)_+^2 \right]} \\ & \text{subject to} \quad \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2} \leq \alpha, \quad \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} = \theta \end{aligned} \quad (\text{B.21})$$

Bringing the constraints into the objective via Lagrange multipliers we obtain

$$\begin{aligned} & \min_{\boldsymbol{\theta}, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \lambda, \eta \geq 0, \tilde{\eta}} \quad \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon \gamma - \theta \|\boldsymbol{\mu}\|_{\ell_2}) + \alpha g \right)_+^2 \right]} \\ & + \eta \left(\left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2} - \alpha \right) + \tilde{\eta} (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} - \theta) \end{aligned} \quad (\text{B.22})$$

Next note that $\left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2} = \min_{\tau \geq 0} \frac{\left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2}^2}{2\tau} + \frac{\tau}{2}$ and $\|\boldsymbol{\theta}\|_{\ell_q} = \min_{t \geq 0} \frac{\|\boldsymbol{\theta}\|_{\ell_q}^q}{qt^{q-1}} + \frac{q-1}{q}t$

Thus, above reduces to

$$\begin{aligned} & \min_{\boldsymbol{\theta}, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \lambda, \eta \geq 0, \tilde{\eta}} \min_{\tau \geq 0, t \geq 0} \quad \|\boldsymbol{\theta}\|_{\ell_2}^2 + \frac{2\lambda}{qt^{q-1}} \|\boldsymbol{\theta}\|_{\ell_q}^q + 2\lambda \frac{q-1}{q}t - 2\lambda\gamma + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \\ & + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon \gamma - \theta \|\boldsymbol{\mu}\|_{\ell_2}) + \alpha g \right)_+^2 \right]} \\ & + \frac{\eta}{2\tau} \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2}^2 + \frac{\eta\tau}{2} - \eta\alpha + \tilde{\eta} (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} - \theta) \end{aligned} \quad (\text{B.23})$$

To continue note that $\frac{\|\boldsymbol{\theta}\|_{\ell_q}^q}{t^{q-1}} = t \left\| \frac{\boldsymbol{\theta}}{t} \right\|_{\ell_q}^q$ and $\frac{\left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2}^2}{\tau} = \tau \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \frac{\boldsymbol{\theta}}{\tau} \right\|_{\ell_2}^2$ and thus using the fact that the perspective of a convex function is convex both are

jointly convex with respect to $(\boldsymbol{\theta}, t)$ and $(\boldsymbol{\theta}, \tau)$. Thus the objective above is jointly convex in $(\boldsymbol{\theta}, \theta, \gamma, \alpha, t, \tau)$ and jointly concave in $(\beta, \lambda, \eta, \tilde{\eta})$. Due to this convexity/concavity with respect to the minimization/maximization parameters we can change the order of min and max. We thus proceed by optimizing over $\boldsymbol{\theta}$. The optimization over $\boldsymbol{\theta}$ takes the form

$$(B.24) \quad \min_{\boldsymbol{\theta}} \quad \|\boldsymbol{\theta}\|_{\ell_2}^2 + \frac{2\lambda}{qt^{q-1}} \|\boldsymbol{\theta}\|_{\ell_q}^q + \frac{\eta}{2\tau} \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2}^2 + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + \tilde{\eta} \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta}$$

By completing the square the objective can be alternatively written as

$$\begin{aligned} & \boldsymbol{\theta}^T \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right) \boldsymbol{\theta} + \frac{2\lambda}{qt^{q-1}} \|\boldsymbol{\theta}\|_{\ell_q}^q + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + \tilde{\eta} \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} \\ &= \left\| \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{1/2} \boldsymbol{\theta} + \frac{\beta}{2\sqrt{n}} \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1/2} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} + \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1/2} \frac{\tilde{\eta}}{2} \tilde{\boldsymbol{\mu}} \right\|_{\ell_2}^2 \\ &+ \frac{2\lambda}{qt^{q-1}} \|\boldsymbol{\theta}\|_{\ell_q}^q - \frac{\beta^2}{4n} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \boldsymbol{\Sigma}^{1/2} \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} \\ (B.25) \quad & - \frac{\tilde{\eta}^2}{4} \tilde{\boldsymbol{\mu}}^T \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1} \tilde{\boldsymbol{\mu}} - \frac{\beta \tilde{\eta}}{2\sqrt{n}} \tilde{\boldsymbol{\mu}}^T \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h}. \end{aligned}$$

Since $\boldsymbol{\Sigma} \tilde{\boldsymbol{\mu}} = a^2 \tilde{\boldsymbol{\mu}}$ we have

$$\tilde{\boldsymbol{\mu}}^T \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} = 0, \quad \tilde{\boldsymbol{\mu}}^T \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1} \tilde{\boldsymbol{\mu}} = \frac{1}{(1 + \frac{\eta}{2\tau} a^2)}.$$

We consider a singular value decomposition $\boldsymbol{\Sigma} = \mathbf{U} \mathbf{S} \mathbf{U}^T$ with $\mathbf{S} = \text{diag}(s_1, \dots, s_d)$, and the first column of \mathbf{U} being $\tilde{\boldsymbol{\mu}}$ and $s_1 = a^2$ (Recall that $\tilde{\boldsymbol{\mu}}$ is a singular value of $\boldsymbol{\Sigma}$ with eigenvalue a^2 .) Then,

$$\begin{aligned} \frac{1}{n} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \boldsymbol{\Sigma}^{1/2} \left(\mathbf{I} + \frac{\eta}{2\tau} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} &= \frac{1}{n} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \mathbf{U} \left(\mathbf{I} + \frac{\eta}{2\tau} \mathbf{S} \right)^{-1} \mathbf{S} \mathbf{U}^T \mathbf{P}_{\mu}^{\perp} \mathbf{h} \\ &= \frac{1}{\delta d} \sum_{i=2}^d \frac{s_i}{1 + \frac{\eta}{2\tau} s_i} h_i^2 \\ &\stackrel{P}{\Rightarrow} \frac{1}{\delta d} \sum_{i=1}^d \frac{s_i}{1 + \frac{\eta}{2\tau} s_i} \\ &= \frac{2\tau}{\delta d \eta} \sum_{i=1}^d \left(1 - \frac{1}{\frac{\eta}{2\tau} (s_i + \frac{2\tau}{\eta})} \right) \\ &= \frac{2\tau}{\delta \eta} \left(1 + \frac{2\tau}{\eta} S_{\rho} \left(-\frac{2\tau}{\eta} \right) \right) \end{aligned}$$

with $S_\rho(z) := \int \frac{\rho(t)}{z-t} dt$ the Stieltjes transform of the spectrum of Σ .

Using the above identities (B.25) reduces to

$$\begin{aligned} & \left\| \left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{1/2} \boldsymbol{\theta} + \frac{\beta}{2\sqrt{n}} \left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{-1/2} \Sigma^{1/2} \mathbf{P}_\mu^\perp \mathbf{h} + \left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{-1/2} \frac{\tilde{\eta}}{2} \tilde{\boldsymbol{\mu}} \right\|_{\ell_2}^2 \\ & + \frac{2\lambda}{qt^{q-1}} \|\boldsymbol{\theta}\|_{\ell_q}^q - \frac{\tau\beta^2}{2\delta\eta} \left(1 + \frac{2\tau}{\eta} S_\rho \left(-\frac{2\tau}{\eta} \right) \right) - \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\tau} a^2)}. \end{aligned}$$

We then write the minimum value over $\boldsymbol{\theta}$ in terms of the weighted Moreau envelope, given by Definition 6.2.

$$\begin{aligned} & \min_{\boldsymbol{\theta}} \left\| \left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{1/2} \boldsymbol{\theta} + \frac{\beta}{2\sqrt{n}} \left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{-1/2} \Sigma^{1/2} \mathbf{P}_\mu^\perp \mathbf{h} + \left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{-1/2} \frac{\tilde{\eta}}{2} \tilde{\boldsymbol{\mu}} \right\|_{\ell_2}^2 + \frac{2\lambda}{qt^{q-1}} \|\boldsymbol{\theta}\|_{\ell_q}^q \\ & \text{(B.26)} \\ & = 2e_{q, \mathbf{I} + \frac{\eta}{2\tau} \Sigma} \left(\left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{-1} \left\{ \frac{\beta}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_\mu^\perp \mathbf{h} - \frac{\tilde{\eta}}{2} \tilde{\boldsymbol{\mu}} \right\}; \frac{\lambda}{qt^{q-1}} \right), \end{aligned}$$

where we used symmetry of the distribution of \mathbf{h} .

Putting all pieces together in (B.25) we get

$$\begin{aligned} & \text{(B.27)} \\ & \min_{\boldsymbol{\theta}} \|\boldsymbol{\theta}\|_{\ell_2}^2 + \frac{2\lambda}{qt^{q-1}} \|\boldsymbol{\theta}\|_{\ell_q}^q + \frac{\eta}{2\tau} \left\| \Sigma^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2}^2 + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_\mu^\perp \Sigma^{1/2} \boldsymbol{\theta} + \tilde{\eta} \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} \\ & = 2e_{q, \mathbf{I} + \frac{\eta}{2\tau} \Sigma} \left(\left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{-1} \left\{ \frac{\beta}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_\mu^\perp \mathbf{h} - \frac{\tilde{\eta}}{2} \tilde{\boldsymbol{\mu}} \right\}; \frac{\lambda}{qt^{q-1}} \right) - \frac{\beta^2 \tau}{2\delta\eta} \left(1 + \frac{2\tau}{\eta} S_\rho \left(-\frac{2\tau}{\eta} \right) \right) \\ & - \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\tau} a^2)}. \end{aligned}$$

Using (B.27) in (B.23), the AO problem reduces to

$$\begin{aligned} & \min_{\gamma \geq 0, \theta} \max_{\beta, \lambda, \eta \geq 0, \tilde{\eta}} \min_{\tau \geq 0, t \geq 0} 2e_{q, \mathbf{I} + \frac{\eta}{2\tau} \Sigma} \left(\left(\mathbf{I} + \frac{\eta}{2\tau} \Sigma \right)^{-1} \left\{ \frac{\beta}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_\mu^\perp \mathbf{h} - \frac{\tilde{\eta}}{2} \tilde{\boldsymbol{\mu}} \right\}; \frac{\lambda}{qt^{q-1}} \right) \\ & - \frac{\beta^2 \tau}{2\delta\eta} \left(1 + \frac{2\tau}{\eta} S_\rho \left(-\frac{2\tau}{\eta} \right) \right) \\ & - \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\tau} a^2)} + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon\gamma - \theta \|\boldsymbol{\mu}\|_{\ell_2}) + \alpha g \right)_+^2 \right]} \\ & \text{(B.28)} \\ & + 2\lambda \frac{q-1}{q} t - 2\lambda\gamma + \frac{\eta\tau}{2} - \eta\alpha - \tilde{\eta}\theta \end{aligned}$$

Scalarization of the auxiliary optimization problem. We proceed by defining $\lambda_0 := \frac{\lambda}{\|\boldsymbol{\mu}\|_{\ell_p}}$, $\gamma_0 := \gamma \|\boldsymbol{\mu}\|_{\ell_p}$ and $t_0 := t \|\boldsymbol{\mu}\|_{\ell_p}$. Under Assumptions 2 and 6, the asymptotic auxiliary optimization (AO) problem becomes

$$\begin{aligned}
 (B.29) \quad & \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\beta, \lambda_0, \eta \geq 0, \tilde{\eta}} \min_{\tau \geq 0, t_0 \geq 0} 2F\left(\beta, \tilde{\eta}; \frac{\eta}{2\tau}, \frac{\lambda_0}{qt_0^{q-1}}\right) - \frac{\beta^2 \tau}{2\delta\eta} \left(1 + \frac{2\tau}{\eta} S_\rho\left(-\frac{2\tau}{\eta}\right)\right) \\
 & + 2\lambda_0 \frac{q-1}{q} t_0 - 2\lambda_0 \gamma_0 + \frac{\eta\tau}{2} - \eta\alpha - \tilde{\eta}\theta \\
 & - \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\tau} a^2)} + \beta \sqrt{\mathbb{E}\left[\left((1 + \varepsilon_0 \gamma_0 - \theta V) + \alpha g\right)_+^2\right]}
 \end{aligned}$$

Here we used the relation $\varepsilon = \varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}$.

We next solve for some of the variables in the AO problem by writing the KKT conditions.

1. Define

$$f\left(\frac{\eta}{\tau}\right) := 2F\left(\beta, \tilde{\eta}; \frac{\eta}{2\tau}, \frac{\lambda_0}{qt_0^{q-1}}\right) - \frac{\beta^2 \tau}{2\delta\eta} \left(1 + \frac{2\tau}{\eta} S_\rho\left(-\frac{2\tau}{\eta}\right)\right) - \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\tau} a^2)},$$

where we only made the dependence on $\frac{\eta}{\tau}$ explicit in the notation $f\left(\frac{\eta}{\tau}\right)$. Setting derivative with respect to η to zero, we obtain

$$(B.30) \quad \frac{1}{\tau} f'\left(\frac{\eta}{\tau}\right) + \frac{\tau}{2} - \alpha = 0.$$

Setting derivative with respect to τ to zero, we obtain

$$(B.31) \quad -\frac{\eta}{\tau^2} f'\left(\frac{\eta}{\tau}\right) + \frac{\eta}{2} = 0.$$

Combining (B.30) and (B.31), we get $\eta(1 - \frac{\alpha}{\tau}) = 0$. So either $\alpha = \tau$ or $\eta = 0$. If $\eta = 0$, then it is clear that the terms involving τ in the AO problem would vanish and therefore the value of τ does not matter. So in this case, we can as well assume $\tau = \alpha$. This simplifies the AO problem by replacing for τ :

$$\begin{aligned}
 (B.32) \quad & \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\beta, \lambda_0, \eta \geq 0, \tilde{\eta}} \min_{t_0 \geq 0} 2F\left(\beta, \tilde{\eta}; \frac{\eta}{2\alpha}, \frac{\lambda_0}{qt_0^{q-1}}\right) - \frac{\beta^2 \alpha}{2\delta\eta} \left(1 + \frac{2\alpha}{\eta} S_\rho\left(-\frac{2\alpha}{\eta}\right)\right) \\
 & + 2\lambda_0 \frac{q-1}{q} t_0 - 2\lambda_0 \gamma_0 - \frac{\eta\alpha}{2} - \tilde{\eta}\theta \\
 & - \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\alpha} a^2)} + \beta \sqrt{\mathbb{E}\left[\left((1 + \varepsilon_0 \gamma_0 - \theta V) + \alpha g\right)_+^2\right]}.
 \end{aligned}$$

2. Setting derivative with respect to λ_0 to zero, we get

$$(B.33) \quad 2F'_4\left(\beta, \tilde{\eta}; \frac{\eta}{2\alpha}, \frac{\lambda_0}{qt_0^{q-1}}\right) \frac{1}{qt_0^{q-1}} + 2\frac{q-1}{q}t_0 - 2\gamma_0 = 0,$$

where F'_4 denotes the derivative of function F with respect to its forth argument. Also, by setting derivative with respect to t_0 to zero we get

$$(B.34) \quad 2F'_4\left(\beta, \tilde{\eta}; \frac{\eta}{2\alpha}, \frac{\lambda_0}{qt_0^{q-1}}\right) \lambda_0 \frac{1-q}{q} t_0^{-q} + 2\lambda_0 \frac{q-1}{q} = 0.$$

Combining (B.33) and (B.34) implies that

$$(B.35) \quad 2\lambda_0(q-1)\left(\frac{\gamma_0}{t_0} - 1\right) = 0.$$

Therefore either $\gamma_0 = t_0$ or $\lambda_0 = 0$ or $q = 1$. If $\lambda = 0$ or $q = 1$ then the terms involving t_0 in (B.32) vanish and hence we can assume $t_0 = \gamma_0$ in this cases as well. Replacing t_0 with γ_0 in (B.32) we obtain

$$(B.36) \quad \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\beta, \lambda_0, \eta \geq 0, \tilde{\eta}} 2F\left(\beta, \tilde{\eta}; \frac{\eta}{2\alpha}, \frac{\lambda_0}{q\gamma_0^{q-1}}\right) - \frac{\beta^2\alpha}{2\delta\eta} \left(1 + \frac{2\alpha}{\eta} S_\rho\left(-\frac{2\alpha}{\eta}\right)\right) - \frac{2\lambda_0}{q}\gamma_0 - \frac{\eta\alpha}{2} - \tilde{\eta}\theta$$

$$- \frac{\tilde{\eta}^2}{4(1 + \frac{\eta}{2\alpha}a^2)} + \beta \sqrt{\mathbb{E}\left[\left((1 + \varepsilon_0\gamma_0 - \theta V) + \alpha g\right)_+^2\right]},$$

which is the expression for $D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta})$ given by (6.2).

Uniqueness and boundedness of the solution to AO problem. Note that since $\delta \leq \delta_*$, by using Theorem 6.1, we are in the separable regime and therefore optimization (F.2) is feasible with high probability and admits a bounded solution. This implies that the PO problem (B.17) has bounded solution and since AO and PO problems are asymptotically equivalent this implies that the AO problem (B.36) has bounded solution.

To show the uniqueness of the solution of (B.36), note that as we argued throughout the proof, its objective function D_s is jointly strictly convex in $(\alpha, \gamma_0, \theta)$ and jointly concave in $(\beta, \lambda_0, \eta, \tilde{\eta})$. Therefore, $\max_{\beta, \lambda_0, \eta, \tilde{\eta}} D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta})$ is strictly convex in $(\alpha, \gamma_0, \theta)$. This follows from the fact that if a function $f(\mathbf{x}, \mathbf{y})$ is strictly convex in \mathbf{x} , then $\max_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is also strictly convex in \mathbf{x} and therefore has a unique minimizer $(\alpha_*, \gamma_{0*}, \theta_*)$.

Part (b) of the theorem follows readily from our definition of parameters α , θ and γ .

Part (c) also follows from combining Lemma 2.1 with part (b) of the theorem.

B.3. Proof of Theorem 6.4 The goal of this theorem is to derive precise asymptotic behavior for the adversarially trained model $\widehat{\theta}^\varepsilon$ given by

$$(B.37) \quad \widehat{\theta}^\varepsilon = \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell \left(y_i \langle \mathbf{x}_i, \theta \rangle - \varepsilon \|\theta\|_{\ell_q} \right).$$

Letting $v_i := y_i \langle \mathbf{x}_i, \theta \rangle$, this optimization can be equivalently written as

$$\min_{\theta, v \in \mathbb{R}^n} \frac{1}{2p} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\theta\|_{\ell_q} \right) \quad \text{subject to } v = D_y X \theta,$$

with $D_y = \text{diag}(y_1, \dots, y_n)$. Therefore, by writing the Lagrangian by \mathbf{u}/n as the dual variable for the equality constraint, we arrive at

$$\min_{\theta, v \in \mathbb{R}^n} \max_{\mathbf{u} \in \mathbb{R}^n} \frac{1}{n} \left\{ \mathbf{u}^\top D_y X \theta - \mathbf{u}^\top v \right\} + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\theta\|_{\ell_q} \right)$$

We next substitute for $X = \mathbf{y} \mu^\top + Z \Sigma^{1/2}$, under the Gaussian mixtures model, which gives us

$$(B.38) \quad \min_{\theta, v \in \mathbb{R}^n} \max_{\mathbf{u} \in \mathbb{R}^n} \frac{1}{n} \left\{ \mathbf{u}^\top \mathbf{1} \mu^\top \theta + \mathbf{u}^\top D_y Z \Sigma^{1/2} \theta - \mathbf{u}^\top v \right\} + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\theta\|_{\ell_q} \right)$$

Note that by the above Lagrangian is in a minimax problem in the form of $\min_{\theta} \max_{\mathbf{u}} \mathbf{u}^\top Z \theta + \psi(\theta, \mathbf{u})$, with Z standard Gaussian matrix and $\psi(\theta, \mathbf{u})$ is convex in the minimization variable θ and concave in the maximization variable \mathbf{u} . This form allows us to apply the CGMT framework as outlined in Section 7 and similar to the proof of Theorem 6.3. But in order to do that, we need the minimization/maximization to be over compact sets. Similar to the proof of Theorem 6.3 we cope with this technical issue by introducing artificial boundedness constraints on the optimization variables that they do not change the optimal solution. Specifically, we can add constraints of the form $\mathcal{S}_\theta = \{\theta : \|\theta\|_{\ell_q} \leq K_\theta\}$ and $\mathcal{S}_\mathbf{u} = \{\mathbf{u} : \|\mathbf{u}\|_\infty \leq K_\mathbf{u}\}$ for sufficiently large constants $K_\theta, K_\mathbf{u}$, without changing the optimal solution of (B.38). We refer to Appendix E.3.1 for precise statements and proofs. This allows us to

replace (B.38) with

(B.39)

$$\min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v} \in \mathbb{R}^n} \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} \frac{1}{n} \left\{ \mathbf{u}^\top \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} + \mathbf{u}^\top \mathbf{D}_{\mathbf{y}} \mathbf{Z} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} - \mathbf{u}^\top \mathbf{v} \right\} + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right).$$

B.3.1. *Reduction to an auxiliary optimization problem via CGMT* Next we define the projection matrices

$$\mathbf{P}_{\boldsymbol{\mu}}^\perp := \mathbf{I} - \tilde{\boldsymbol{\mu}} \tilde{\boldsymbol{\mu}}^T, \quad \mathbf{P}_{\boldsymbol{\mu}} := \tilde{\boldsymbol{\mu}} \tilde{\boldsymbol{\mu}}^T$$

and rewrite $\mathbf{Z} \boldsymbol{\Sigma}^{1/2} = \mathbf{Z} (\mathbf{P}_{\boldsymbol{\mu}} + \mathbf{P}_{\boldsymbol{\mu}}^\perp) \boldsymbol{\Sigma}^{1/2}$. Since $\mathbf{Z} \mathbf{P}_{\boldsymbol{\mu}}$ and $\mathbf{Z} \mathbf{P}_{\boldsymbol{\mu}}^\perp$ are independent from each other the latter has the same distribution as

$$(B.40) \quad \mathbf{Z} \boldsymbol{\Sigma}^{1/2} \sim \mathbf{z} \left(\boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} \right)^T + \mathbf{Z} \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2}.$$

where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_n)$ and is independent from the matrix \mathbf{Z} . Thus the above optimization problem is equivalent to

(B.41)

$$\min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v} \in \mathcal{S}_{\mathbf{u}}} \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} \frac{1}{n} \left\{ \mathbf{u}^\top \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} + (\mathbf{u}^\top \mathbf{D}_{\mathbf{y}} \mathbf{z}) (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) + \mathbf{u}^\top \mathbf{D}_{\mathbf{y}} \mathbf{Z} \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} - \mathbf{u}^\top \mathbf{v} \right\} + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right).$$

Using CGMT and the corresponding AO takes the form

$$(B.42) \quad \min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v} \in \mathcal{S}_{\mathbf{u}}} \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} \frac{1}{n} \left\{ \left\| \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g}^T \mathbf{D}_{\mathbf{y}} \mathbf{u} + \left\| \mathbf{D}_{\mathbf{y}} \mathbf{u} \right\|_{\ell_2} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + (\mathbf{u}^\top \mathbf{D}_{\mathbf{y}} \mathbf{z}) (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) \right. \\ \left. + \mathbf{u}^\top \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} - \mathbf{u}^\top \mathbf{v} \right\} + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right),$$

where $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_n)$ and $\mathbf{h} \sim \mathcal{N}(0, \mathbf{I}_d)$.

Given $y_i = \pm 1$ are independent of \mathbf{Z} and hence \mathbf{g} , we have $\mathbf{D}_{\mathbf{y}} \mathbf{g}, \mathbf{D}_{\mathbf{y}} \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_n)$ and $\left\| \mathbf{D}_{\mathbf{y}} \mathbf{u} \right\|_{\ell_2} = \left\| \mathbf{u} \right\|_{\ell_2}$. This results in

$$(B.43) \quad \min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v} \in \mathcal{S}_{\mathbf{u}}} \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} \frac{1}{n} \left\{ \left\| \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g}^T \mathbf{u} + \left\| \mathbf{u} \right\|_{\ell_2} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} + (\mathbf{u}^\top \mathbf{z}) (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}) + \mathbf{u}^\top \mathbf{1} \boldsymbol{\mu}^T \boldsymbol{\theta} - \mathbf{u}^\top \mathbf{v} \right\} \\ + \frac{1}{n} \sum_{i=1}^n \ell \left(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right).$$

Letting $\beta := \frac{1}{\sqrt{n}} \|\mathbf{u}\|_{\ell_2}$ and optimizing over direction of \mathbf{u} , we get

$$\begin{aligned} & \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} \frac{1}{n} \left(\left\| \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g}^T \mathbf{u} + \|\mathbf{u}\|_{\ell_2} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} + (\mathbf{u}^T \mathbf{z})(\tilde{\mu}^T \Sigma^{1/2} \boldsymbol{\theta}) + \mathbf{u}^T \mathbf{1} \mu^T \boldsymbol{\theta} - \mathbf{u}^T \mathbf{v} \right) \\ & \quad (\text{B.44}) \\ & = \max_{0 \leq \beta \leq K} \frac{\beta}{\sqrt{n}} \left\| \left\| \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} + \tilde{\mu}^T \Sigma^{1/2} \boldsymbol{\theta} \mathbf{z} + \mathbf{1} \mu^T \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2} + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta}, \end{aligned}$$

where $K := \max_{\mathbf{u} \in \mathcal{S}_{\mathbf{u}}} \frac{1}{\sqrt{n}} \|\mathbf{u}\|_{\ell_2} < K_{\mathbf{u}}$ by definition of $\mathcal{S}_{\mathbf{u}}$.

Plugging the latter into AO becomes

$$\begin{aligned} & \quad (\text{B.45}) \\ & \min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v}} \max_{0 \leq \beta \leq K} \frac{\beta}{\sqrt{n}} \left\| \left\| \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} + \tilde{\mu}^T \Sigma^{1/2} \boldsymbol{\theta} \mathbf{z} + \mathbf{1} \mu^T \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2} + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} + \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q}). \end{aligned}$$

We hereafter use the shorthand

$$\ell(\mathbf{v}, \boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q}),$$

for simplicity of notation. For the minimization, with respect to $\boldsymbol{\theta}$ and then \mathbf{v} , to become easier in our later calculation we proceed by writing $\ell(\mathbf{v}, \boldsymbol{\theta})$ in terms of its conjugate with respect to $\boldsymbol{\theta}$. That is,

$$\ell(\mathbf{v}, \boldsymbol{\theta}) = \sup_{\mathbf{w}} \mathbf{w}^T \boldsymbol{\theta} - \tilde{\ell}(\mathbf{v}, \mathbf{w})$$

where $\tilde{\ell}(\mathbf{v}, \mathbf{w})$ is the conjugate of ℓ with respect to $\boldsymbol{\theta}$. The logic behind this is that AO will then simplify to

$$\min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v}} \max_{0 \leq \beta \leq K, \mathbf{w}} \frac{\beta}{\sqrt{n}} \left\| \left\| \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} + \tilde{\mu}^T \Sigma^{1/2} \boldsymbol{\theta} \mathbf{z} + \mathbf{1} \mu^T \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2} + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} + \mathbf{w}^T \boldsymbol{\theta} - \tilde{\ell}(\mathbf{v}, \mathbf{w})$$

which after flipping (allowed based on the correct form of convexity/concavity of PO) becomes

$$\begin{aligned} & \quad (\text{B.46}) \\ & \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v}} \frac{\beta}{\sqrt{n}} \left\| \left\| \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \mathbf{g} + \tilde{\mu}^T \Sigma^{1/2} \boldsymbol{\theta} \mathbf{z} + \mathbf{1} \mu^T \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2} + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\mu}^{\perp} \Sigma^{1/2} \boldsymbol{\theta} + \mathbf{w}^T \boldsymbol{\theta} - \tilde{\ell}(\mathbf{v}, \mathbf{w}). \end{aligned}$$

We define the parallel and perpendicular components of $\boldsymbol{\theta}$ along vector μ as follows:

$$\begin{aligned} & \quad (\text{B.47}) \\ & \boldsymbol{\theta}_{\perp} = \mathbf{P}_{\mu}^{\perp} \boldsymbol{\theta}, \quad \boldsymbol{\theta} := \tilde{\mu}^T \boldsymbol{\theta}, \quad \mathbf{P}_{\mu} \boldsymbol{\theta} = \boldsymbol{\theta} \tilde{\mu}. \end{aligned}$$

Given that $\tilde{\boldsymbol{\mu}}$ is an eigenvector of $\boldsymbol{\Sigma}$, cf. Assumption 5, we have $\mathbf{P}_\mu \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp = 0$ and therefore

$$\mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} = \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} (\mathbf{P}_\mu + \mathbf{P}_\mu^\perp) \boldsymbol{\theta} = \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp \boldsymbol{\theta} = \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp \boldsymbol{\theta}_\perp.$$

Similarly, since $\boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} = a \tilde{\boldsymbol{\mu}}$.

$$\tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} = a \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} = a \theta.$$

Rewriting the AO problem, we get

$$(B.48) \quad \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v}} \frac{\beta}{\sqrt{n}} \left\| \left\| \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp \boldsymbol{\theta}_\perp \right\|_{\ell_2} \mathbf{g} + a \theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2} \\ + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp \boldsymbol{\theta}_\perp + \mathbf{w}^T \mathbf{P}_\mu^\perp \boldsymbol{\theta}_\perp + \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta - \tilde{\ell}(\mathbf{v}, \mathbf{w}).$$

We can rewrite this as

$$(B.49) \quad \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}, \mathbf{v}} \frac{\beta}{\sqrt{n}} \left\| \left\| \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp \boldsymbol{\theta}_\perp \right\|_{\ell_2} \mathbf{g} + a \theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2} \\ + \frac{\beta}{\sqrt{n}} (\mathbf{P}_\mu^\perp \mathbf{h})^T \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp \boldsymbol{\theta}_\perp + \mathbf{w}^T \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_\mu^\perp (\mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp) \boldsymbol{\theta}_\perp + \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta - \tilde{\ell}(\mathbf{v}, \mathbf{w}).$$

Here we used the assumption that $\tilde{\boldsymbol{\mu}}$ is an eigenvector of $\boldsymbol{\Sigma}$ which in turn implies that

$$(B.50) \quad \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_\mu^\perp (\mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp) = \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp = \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{-1/2} (\mathbf{P}_\mu + \mathbf{P}_\mu^\perp) \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp = \mathbf{P}_\mu^\perp.$$

We next optimize over $\boldsymbol{\theta}$ using lemma below and its proof is deferred to Appendix E.4.

Lemma B.1 *For a given vector \mathbf{r} and $\alpha \geq 0$ consider the following optimization*

$$(B.51) \quad \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \langle \mathbf{P}_\mu^\perp \mathbf{r}, (\mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp) \boldsymbol{\theta}_\perp \rangle$$

$$(B.52) \quad \text{subject to } \left\| (\mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp) \boldsymbol{\theta}_\perp \right\|_{\ell_2} = \alpha$$

Under the assumption that $\mathbf{P}_\mu \boldsymbol{\Sigma}^{1/2} \mathbf{P}_\mu^\perp = 0$, the optimal value of this optimization is given by $-\alpha \|\mathbf{P}_\mu^\perp \mathbf{r}\|_{\ell_2}$.

Now note that

$$|\theta| = |\tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta}| \leq \|\boldsymbol{\theta}\|_{\ell_2},$$

$$\alpha = \left\| \left(\mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\boldsymbol{\mu}}^\perp \right) \boldsymbol{\theta} \right\|_{\ell_2} = \left\| \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \leq \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \right\|_{\ell_2} \leq C_{\max}^{1/2} \|\boldsymbol{\theta}\|_{\ell_2}$$

where in the second line we used Assumption 1(b),(d). Since $\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}$ a bounded set, we can choose $K' > 0$ large enough so that $0 \leq |\theta|, \alpha \leq K'$ and hence so do the optimization over this bounded range. That said, we use Lemma B.1 with $\mathbf{r} = \frac{\beta}{\sqrt{n}} \mathbf{h} + \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^\perp \mathbf{w}$, to simplify the AO problem as follows:

$$(B.53) \quad \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{0 \leq \alpha, |\theta| \leq K', \mathbf{v}} \frac{\beta}{\sqrt{n}} \left\| \alpha \mathbf{g} + a\theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2}$$

$$- \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^\perp \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^\perp \mathbf{w} \right\|_{\ell_2} + \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta - \tilde{\ell}(\mathbf{v}, \mathbf{w})$$

To continue we shall calculate the conjugate function $\tilde{\ell}$. This is the subject of the next lemma and we refer to Appendix E.5 for its proof.

Lemma B.2 *The conjugate of the function*

$$\ell(\mathbf{v}, \boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q})$$

with respect to $\boldsymbol{\theta}$ is equal to

$$\tilde{\ell}(\mathbf{v}, \mathbf{w}) = \sup_{\gamma \geq 0} \gamma \|\mathbf{w}\|_{\ell_p} - \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \gamma).$$

Using the above lemma we have

$$-\tilde{\ell}(\mathbf{v}, \mathbf{w}) = - \left(\sup_{\gamma \geq 0} \gamma \|\mathbf{w}\|_{\ell_p} - \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \gamma) \right) = \inf_{\gamma \geq 0} -\gamma \|\mathbf{w}\|_{\ell_p} + \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \gamma).$$

Plugging this into (B.53) we arrive at

$$(B.54) \quad \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{0 \leq \alpha, |\theta| \leq K', \mathbf{v}, 0 \leq \gamma} \frac{\beta}{\sqrt{n}} \left\| \alpha \mathbf{g} + a\theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \boldsymbol{\theta} - \mathbf{v} \right\|_{\ell_2}$$

$$- \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^\perp \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^\perp \mathbf{w} \right\|_{\ell_2} + \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta - \gamma \|\mathbf{w}\|_{\ell_p} + \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \gamma)$$

Note that when $p \geq 1$ the objective is jointly concave in (\mathbf{w}, β) and jointly convex in $\alpha, \theta, \mathbf{v}$ and therefore we can switch the orders of min and max.

We next focus on optimization over \mathbf{v} . Using the observation that for all $x \in \mathbb{R}$, $\min_{\tau \geq 0} \frac{\tau}{2} + \frac{x^2}{2\tau} = x$ we write

$$\begin{aligned}
& \min_{\mathbf{v}} \frac{\beta}{\sqrt{n}} \left\| \alpha \mathbf{g} + a\theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \theta - \mathbf{v} \right\|_{\ell_2} + \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon\gamma) \\
&= \min_{\mathbf{v}} \inf_{\tau_g \geq 0} \frac{\beta}{2\tau_g n} \left\| \alpha \mathbf{g} + a\theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \theta - \mathbf{v} \right\|_{\ell_2}^2 + \frac{\beta\tau_g}{2} + \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon\gamma) \\
&= \min_{\mathbf{v}} \inf_{\tau_g \geq 0} \frac{\beta}{2\tau_g n} \sum_{i=1}^n \left(\alpha g_i + a\theta z_i + \|\boldsymbol{\mu}\|_{\ell_2} \theta - v_i \right)^2 + \frac{\beta\tau_g}{2} + \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon\gamma) \\
&\quad (\text{B.55}) \\
&= \min_{\tilde{\mathbf{v}}} \inf_{\tau_g \geq 0} \frac{\beta}{2\tau_g n} \sum_{i=1}^n \left(\alpha g_i + a\theta z_i + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \tilde{v}_i - \varepsilon\gamma \right)^2 + \frac{\beta\tau_g}{2} + \frac{1}{n} \sum_{i=1}^n \ell(\tilde{v}_i)
\end{aligned}$$

As a result (B.54) can be rewritten as

$$\begin{aligned}
& \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{0 \leq \alpha, |\theta| \leq K', \tilde{\mathbf{v}}, 0 \leq \gamma} \inf_{\tau_g \geq 0} \frac{\beta}{2\tau_g n} \left\| \alpha \mathbf{g} + a\theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \theta - \tilde{\mathbf{v}} - \varepsilon\gamma \mathbf{1} \right\|_{\ell_2}^2 + \frac{\beta\tau_g}{2} \\
&\quad (\text{B.56}) \quad - \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2} + \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta - \gamma \|\mathbf{w}\|_{\ell_p} + \frac{1}{n} \sum_{i=1}^n \ell(\tilde{v}_i)
\end{aligned}$$

We note that since the quadratic over linear function is jointly convex the above loss is jointly convex in the parameters $(\alpha, \gamma, \theta, \tau_g, \tilde{\mathbf{v}})$. Also for $p \geq 1$ the $\|\cdot\|_{\ell_p}$ is convex and thus the objective is also jointly concave in (β, \mathbf{w}) .

We recall the definition of the Moreau envelope of function ℓ at a point x with parameter μ , that is given by

$$(\text{B.57}) \quad e_{\ell}(x; \mu) := \min_t \frac{1}{2\mu} (x - t)^2 + \ell(t).$$

We can now rewrite equation (B.56) in terms of Moreau envelope of the loss function ℓ .

$$\begin{aligned}
& \min_{\tilde{\mathbf{v}}} \inf_{\tau_g \geq 0} \frac{\beta}{2\tau_g n} \left\| \alpha \mathbf{g} + a\theta \mathbf{z} + \mathbf{1} \|\boldsymbol{\mu}\|_{\ell_2} \theta - \tilde{\mathbf{v}} - \varepsilon\gamma \mathbf{1} \right\|_{\ell_2}^2 + \frac{\beta\tau_g}{2} + \frac{1}{n} \sum_{i=1}^n \ell(\tilde{v}_i) \\
&\quad (\text{B.58}) \quad = \inf_{\tau_g \geq 0} \frac{\beta\tau_g}{2} + \frac{1}{n} \sum_{i=1}^n e_{\ell} \left(\alpha g_i + a\theta z_i + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma; \frac{\tau_g}{\beta} \right)
\end{aligned}$$

Thus (B.56) can be rewritten in the form

$$\begin{aligned}
& \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{0 \leq \alpha, |\theta| \leq K', 0 \leq \gamma} \inf_{\tau_g \geq 0} \frac{\beta\tau_g}{2} + \frac{1}{n} \sum_{i=1}^n e_{\ell} \left(\alpha g_i + a\theta z_i + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma; \frac{\tau_g}{\beta} \right) \\
&\quad (\text{B.59}) \quad - \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2} + \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta - \gamma \|\mathbf{w}\|_{\ell_p}
\end{aligned}$$

Note that since (B.56) is jointly convex in $(\alpha, \gamma, \theta, \tau_g, \tilde{\mathbf{v}})$ and jointly concave in (β, \mathbf{w}) and partial minimization preserves convexity thus (B.59) is jointly convex in $(\alpha, \gamma, \theta, \tau_g)$ and jointly concave in (β, \mathbf{w}) .

B.3.2. Scalarization of the auxiliary optimization problem The auxiliary problem (B.59) is in terms of high-dimensional vectors $\mathbf{g}, \mathbf{z}, \mathbf{h}, \mathbf{w}, \boldsymbol{\mu}$. We turn this problem into a scalar optimization by taking the pointwise limit of its objective and then showing that such convergence indeed holds in a uniform sense and therefore the minimax value also converges to that of the limit objective.

Note that by definition of the Moreau envelope, for all x and μ we have

$$e_\ell(x; \mu) \leq \frac{1}{2\mu}(x - x)^2 + \ell(x) = \ell(x) = \log(1 + e^{-x}) \leq \log 2 + |x|.$$

Hence,

$$\mathbb{E} \left[e_\ell \left(\alpha g + a\theta z + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma; \frac{\tau_g}{\beta} \right) \right] \leq \log 2 + \mathbb{E} [|\alpha g + a\theta z + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma|] < \infty,$$

for any finite value of α, θ and γ . Therefore by an application of the Weak Law of Large Numbers, we have that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n e_\ell \left(\alpha g_i + a\theta z_i + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma; \frac{\tau_g}{\beta} \right) &\rightarrow \mathbb{E} \left[e_\ell \left(\alpha g + a\theta z + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma; \frac{\tau_g}{\beta} \right) \right] \\ &= \mathbb{E} \left[e_\ell \left(\sqrt{\alpha^2 + a^2\theta^2} g + \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma; \frac{\tau_g}{\beta} \right) \right]. \end{aligned}$$

We define the expected Moreau envelope $L(a, b, \mu) = \mathbb{E}[e_\ell(ag + b; \mu)]$, where the expectation is taken with respect to independent standard normal variable g .

This simplifies the AO problem as

$$\begin{aligned} \max_{0 \leq \beta \leq K, \mathbf{w}} \min_{0 \leq \alpha, |\theta| \leq K', 0 \leq \gamma, \tau_g} & \frac{\beta\tau_g}{2} + L \left(\sqrt{\alpha^2 + a^2\theta^2}, \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon\gamma, \frac{\tau_g}{\beta} \right) \\ \text{(B.60)} \quad & - \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_\mu^\perp \mathbf{h} + \mathbf{P}_\mu^\perp \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_\mu^\perp \mathbf{w} \right\|_{\ell_2} + \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta - \gamma \|\mathbf{w}\|_{\ell_p}. \end{aligned}$$

We note that since (B.59) is jointly convex in $(\alpha, \gamma, \theta, \tau_g)$ and jointly concave in (β, \mathbf{w}) and expectation preserves convexity thus the objective in (B.60) is also jointly convex in $(\alpha, \gamma, \theta, \tau_g)$ and jointly concave in (β, \mathbf{w}) and

by Sinov's theorem we can flip the maximization over \mathbf{w} and the minimization to arrive at

$$(B.61) \quad \max_{0 \leq \beta \leq K} \min_{0 \leq \alpha, |\theta| \leq K', 0 < \gamma, \tau_g} \frac{\beta \tau_g}{2} + L \left(\sqrt{\alpha^2 + a^2 \theta^2}, \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon \gamma, \frac{\tau_g}{\beta} \right) \\ - \min_{\mathbf{w}} \left\{ \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2} - \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta + \gamma \|\mathbf{w}\|_{\ell_p} \right\}.$$

By our asymptotic setting (cf. Definition 1, part (c)), $\varepsilon = \varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}$ for a constant ε_0 . We let $\gamma_0 := \gamma \|\boldsymbol{\mu}\|_{\ell_p}$ and rewriting (B.61) in terms of γ_0 in lieu of γ we arrive at

$$(B.62) \quad \max_{0 \leq \beta \leq K} \min_{0 \leq \alpha, |\theta| \leq K', 0 < \gamma_0 < K'', 0 < \tau_g} \frac{\beta \tau_g}{2} + L \left(\sqrt{\alpha^2 + a^2 \theta^2}, \|\boldsymbol{\mu}\|_{\ell_2} \theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta} \right) \\ - \min_{\mathbf{w}} \left\{ \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2} - \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta + \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\mathbf{w}\|_{\ell_p} \right\}.$$

• **Optimization over \mathbf{w} .** Continuing with optimization over \mathbf{w} we have

$$(B.63) \quad \min_{\mathbf{w}} \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2} - \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta + \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\mathbf{w}\|_{\ell_p} \\ = \min_{\mathbf{w}, \tau_h \geq 0} \frac{\alpha}{2\tau_h} \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2}^2 + \frac{\alpha \tau_h}{2} - \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta + \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\mathbf{w}\|_{\ell_p} \\ = \min_{\mathbf{w}, \tau_h \geq 0} \frac{\alpha}{2\tau_h} \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{w} \right\|_{\ell_2}^2 + \frac{\alpha \tau_h}{2} - \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta + \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\mathbf{w}\|_{\ell_p},$$

where in the last step we used that $\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}} = 0$, which follows from Assumption 5. Note that the above loss is jointly convex in (\mathbf{w}, τ_h) . So that continuing from (B.61) the overall objective is jointly convex in $(\alpha, \gamma, \theta, \tau_g)$ and jointly concave in $(\beta, \mathbf{w}, \tau_h)$.

Let $\tilde{\mathbf{w}} := \boldsymbol{\Sigma}^{-1/2} \mathbf{w}$. The optimization over \mathbf{w} can be written as

$$(B.64) \quad \min_{\tilde{\mathbf{w}}} \frac{1}{2} \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \tilde{\mathbf{w}} \right\|_{\ell_2}^2 + f(\tilde{\mathbf{w}}),$$

where

$$f(\tilde{\mathbf{w}}) := -\langle \tilde{\mathbf{w}}, \boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} \rangle \frac{\theta \tau_h}{\alpha} + \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\boldsymbol{\Sigma}^{1/2} \tilde{\mathbf{w}}\|_{\ell_p}.$$

Let $\tilde{\mathbf{w}}^*$ be the optimal solution. Then,

$$(B.65) \quad -\mathbf{P}_\mu^\perp \left(\frac{\beta}{\sqrt{n}} \mathbf{h} + \tilde{\mathbf{w}}^* \right) \in \partial f(\tilde{\mathbf{w}}^*).$$

By the conjugate subgradient theorem, this implies that

$$\tilde{\mathbf{w}}^* \in \partial f^* \left(-\mathbf{P}_\mu^\perp \left(\frac{\beta}{\sqrt{n}} \mathbf{h} + \tilde{\mathbf{w}}^* \right) \right).$$

Let $\mathbf{t}^* := \frac{\beta}{\sqrt{n}} \mathbf{h} + \tilde{\mathbf{w}}^*$, then writing the above equation in terms of t ,

$$(B.66) \quad \mathbf{t}^* - \frac{\beta}{\sqrt{n}} \mathbf{h} \in \partial f^* \left(-\mathbf{P}_\mu^\perp \mathbf{t}^* \right).$$

Therefore,

$$(B.67) \quad -\mathbf{P}_\mu^\perp \left(\mathbf{t}^* - \frac{\beta}{\sqrt{n}} \mathbf{h} \right) \in -\mathbf{P}_\mu^\perp \partial f^* \left(-\mathbf{P}_\mu^\perp \mathbf{t}^* \right).$$

Equation (B.67) is equivalent to saying that

$$(B.68) \quad \mathbf{t}^* \in \arg \min_{\mathbf{t}} \frac{1}{2} \left\| \mathbf{P}_\mu^\perp \left(\frac{\beta}{\sqrt{n}} \mathbf{h} - \mathbf{t} \right) \right\|_{\ell_2}^2 + f^* \left(-\mathbf{P}_\mu^\perp \mathbf{t} \right).$$

Lemma B.3 For function $f : \mathbb{R}^p \mapsto \mathbb{R}$ given by

$$f(\tilde{\mathbf{w}}) := -\langle \tilde{\mathbf{w}}, \Sigma^{1/2} \tilde{\boldsymbol{\mu}} \rangle \frac{\theta \tau_h}{\alpha} + \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \left\| \Sigma^{1/2} \tilde{\mathbf{w}} \right\|_{\ell_p},$$

its convex conjugate reads as

$$f^*(\mathbf{u}) = \mathbb{1}_S(\mathbf{u}), \quad S := \left\{ \mathbf{u} : \left\| \Sigma^{-1/2} \mathbf{u} + \frac{\tau_h \theta}{\alpha} \tilde{\boldsymbol{\mu}} \right\|_{\ell_q} \leq \frac{\gamma \tau_h}{\alpha} \right\}, \quad \mathbb{1}_S(\mathbf{u}) = \begin{cases} 0 & \text{if } \mathbf{u} \in S \\ \infty & \text{if } \mathbf{u} \notin S \end{cases}$$

The proof of Lemma B.3 is delegated to Appendix E.6.

Define $\mathcal{B} := \{\boldsymbol{\mu}\}^\perp \cap -S$. Then (B.68) implies that

$$(B.69) \quad \mathbf{P}_\mu^\perp \mathbf{t}^* = \mathbf{P}_{\mathcal{B}} \left(\mathbf{P}_\mu^\perp \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) \right).$$

Lemma B.4 For a convex set \mathcal{S} and $\mathcal{B} := \{\boldsymbol{\mu}\}^\perp \cap \mathcal{S}$, we have $\mathbf{P}_{\mathcal{B}} \mathbf{P}_\mu^\perp = \mathbf{P}_{\mathcal{B}}$.

We refer to Appendix E.7 for the proof of Lemma B.4.
Using Lemma B.4 and (B.69) we obtain

$$\mathbf{P}_\mu^\perp \mathbf{t}^* = \mathbf{P}_\mathcal{B} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right).$$

Recalling definition of \mathbf{t}^* this implies

$$(B.70) \quad \mathbf{P}_\mu^\perp \tilde{\mathbf{w}}^* = \mathbf{P}_\mathcal{B} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) - \frac{\beta}{\sqrt{n}} \mathbf{P}_\mu^\perp \mathbf{h}.$$

Now note that for $p > 1$, $\nabla \|\mathbf{w}\|_{\ell_p} = \frac{1}{\|\mathbf{w}\|_{\ell_p}^{p-1}} [|w_1|^{p-1} \text{sign}(w_1), \dots, |w_p|^{p-1} \text{sign}(w_p)]^\top$.

Therefore in this case $\langle \nabla \|\mathbf{w}\|_{\ell_p}, \mathbf{w} \rangle = \|\mathbf{w}\|_{\ell_p}$. Similarly, for $p = 1$ for any $\mathbf{s} \in \partial \|\mathbf{w}\|_{\ell_p}$ we have $\langle \mathbf{s}, \mathbf{w} \rangle = \|\mathbf{w}\|_{\ell_p}$. Therefore, for all $p \geq 1$ for any $\mathbf{s} \in \partial \|\mathbf{w}\|_{\ell_p}$ we have $\langle \mathbf{s}, \mathbf{w} \rangle = \|\mathbf{w}\|_{\ell_p}$.

Therefore, for the defined function f and any $\mathbf{s} \in \partial f(\tilde{\mathbf{w}})$ there is a vector $\tilde{\mathbf{s}} \in \partial \|\mathbf{x}\|_{\ell_p} \big|_{\mathbf{x}=\Sigma^{1/2}\tilde{\mathbf{w}}}$ such that

$$\begin{aligned} \langle \mathbf{s}, \tilde{\mathbf{w}} \rangle &= \left\langle -\Sigma^{1/2} \tilde{\boldsymbol{\mu}} \frac{\theta \tau_h}{\alpha} + \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \Sigma^{1/2} \tilde{\mathbf{s}}, \tilde{\mathbf{w}} \right\rangle \\ &= -\langle \tilde{\mathbf{w}}, \Sigma^{1/2} \tilde{\boldsymbol{\mu}} \rangle \frac{\theta \tau_h}{\alpha} + \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \langle \tilde{\mathbf{s}}, \Sigma^{1/2} \tilde{\mathbf{w}} \rangle \\ &= -\langle \tilde{\mathbf{w}}, \Sigma^{1/2} \tilde{\boldsymbol{\mu}} \rangle \frac{\theta \tau_h}{\alpha} + \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\Sigma^{1/2} \tilde{\mathbf{w}}\|_{\ell_p} \\ (B.71) \quad &= f(\tilde{\mathbf{w}}). \end{aligned}$$

Therefore by invoking (B.65)

$$\begin{aligned} f(\tilde{\mathbf{w}}^*) &= \left\langle -\mathbf{P}_\mu^\perp \left(\frac{\beta}{\sqrt{n}} \mathbf{h} + \tilde{\mathbf{w}}^* \right), \tilde{\mathbf{w}}^* \right\rangle \\ &= \left\langle -\frac{\beta}{\sqrt{n}} \mathbf{P}_\mu^\perp \mathbf{h} - \mathbf{P}_\mu^\perp \tilde{\mathbf{w}}^*, \mathbf{P}_\mu^\perp \tilde{\mathbf{w}}^* \right\rangle \\ &= \left\langle -\mathbf{P}_\mathcal{B} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right), \mathbf{P}_\mathcal{B} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) - \frac{\beta}{\sqrt{n}} \mathbf{P}_\mu^\perp \mathbf{h} \right\rangle \\ (B.72) \quad &= -\left\| \mathbf{P}_\mathcal{B} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) \right\|_{\ell_2}^2 + \frac{\beta}{\sqrt{n}} \left\langle \mathbf{P}_\mathcal{B} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right), \mathbf{P}_\mu^\perp \mathbf{h} \right\rangle. \end{aligned}$$

Putting things together, the optimal value of objective (B.64) over \mathbf{w} is given by

$$\begin{aligned}
& \min_{\tilde{\mathbf{w}}} \frac{1}{2} \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} + \mathbf{P}_{\mu}^{\perp} \tilde{\mathbf{w}} \right\|_{\ell_2}^2 + f(\tilde{\mathbf{w}}) \\
&= \frac{1}{2} \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} + \mathbf{P}_{\mu}^{\perp} \tilde{\mathbf{w}}^* \right\|_{\ell_2}^2 + f(\tilde{\mathbf{w}}^*) \\
&= \frac{1}{2} \left\| \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) \right\|_{\ell_2}^2 - \left\| \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) \right\|_{\ell_2}^2 + \frac{\beta}{\sqrt{n}} \left\langle \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right), \mathbf{P}_{\mu}^{\perp} \mathbf{h} \right\rangle \\
&= -\frac{1}{2} \left\| \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) \right\|_{\ell_2}^2 + \frac{\beta}{\sqrt{n}} \left\langle \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right), \mathbf{P}_{\mu}^{\perp} \mathbf{h} \right\rangle \\
&= \frac{\beta^2}{2n} \left\| \mathbf{P}_{\mu}^{\perp} \mathbf{h} \right\|_{\ell_2}^2 - \frac{1}{2} \left\| \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) - \frac{\beta}{\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} \right\|_{\ell_2}^2 \\
&= \frac{\beta^2}{2n} \left\| \mathbf{P}_{\mu}^{\perp} \mathbf{h} \right\|_{\ell_2}^2 - \frac{1}{2} \left\| \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) - \frac{\beta}{\sqrt{n}} \mathbf{h} \right\|_{\ell_2}^2 + \frac{\beta^2}{2n} \left\| \mathbf{P}_{\mu} \mathbf{h} \right\|_{\ell_2}^2 \\
\text{(B.73)} \quad &= \frac{\beta^2}{2n} \left\| \mathbf{h} \right\|_{\ell_2}^2 - \frac{1}{2} \left\| \mathbf{P}_{\mathcal{B}} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} \right) - \frac{\beta}{\sqrt{n}} \mathbf{h} \right\|_{\ell_2}^2.
\end{aligned}$$

Following the argument after (B.63) since the objective is jointly convex in $(\alpha, \gamma, \theta, \tau_g)$ and jointly concave in $(\beta, \mathbf{w}, \tau_h)$ and partial maximization preserves concavity after plugging the above the objective is jointly convex in $(\alpha, \gamma, \theta, \tau_g)$ and jointly concave in (β, τ_h) .

- **On projection $\mathbf{P}_{\mathcal{B}}$.** As part of our scalarization process of the auxiliary optimization problem, in the next lemma we provide an alternative characterization of the distance $\|\mathbf{P}_{\mathcal{B}}(\mathbf{h}) - \mathbf{h}\|_{\ell_2}$, and refer to Appendix E.8 for its proof.

Lemma B.5 Recall the set $\mathcal{B} := \{\mu\}^{\perp} \cap -\mathcal{S}$, where \mathcal{S} is given by

$$\mathcal{S} := \left\{ \mathbf{u} : \left\| \Sigma^{-1/2} \mathbf{u} + \frac{\tau_h \theta}{\alpha} \tilde{\mu} \right\|_{\ell_q} \leq \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\mu\|_{\ell_p}} \right\}.$$

Also, suppose that $\Sigma^{1/2} \tilde{\mu} = a \tilde{\mu}$. Then, for any vector \mathbf{h} the following holds:

$$\begin{aligned}
\text{(B.74)} \quad & \frac{1}{2} \left\| \mathbf{P}_{\mathcal{B}}(\mathbf{h}) - \mathbf{h} \right\|_{\ell_2}^2 = \sup_{\lambda \geq 0, \nu} e_{q, \Sigma} \left(\Sigma^{-1/2} \mathbf{h} - \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\mu}; \lambda \right) - \lambda \left(\frac{\gamma_0}{\|\mu\|_{\ell_p}} \frac{\tau_h}{\alpha} \right)^q + \nu \tilde{\mu}^T \mathbf{h} - \frac{\nu^2}{2}
\end{aligned}$$

Using equation (B.73) along with Lemma B.5 we have

$$\begin{aligned}
& \min_{\tilde{\mathbf{w}}} \frac{1}{2} \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \tilde{\mathbf{w}} \right\|_{\ell_2}^2 + f(\tilde{\mathbf{w}}) \\
& \quad (B.75) \\
& = \inf_{\lambda \geq 0, \nu} \frac{\beta^2}{2n} \|\mathbf{h}\|_{\ell_2}^2 + \lambda \left(\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \right)^q - \frac{\nu\beta}{\sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} + \frac{\nu^2}{2} - e_{q, \Sigma} \left(\frac{\beta}{\sqrt{n}} \Sigma^{-1/2} \mathbf{h} - \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\boldsymbol{\mu}}; \lambda \right)
\end{aligned}$$

Recalling equation (B.63) we have

$$\begin{aligned}
& \min_{\mathbf{w}} \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \Sigma^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2} - \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta + \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\mathbf{w}\|_{\ell_p} \\
& = \min_{\tau_h, \lambda \geq 0, \nu} \frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2n} \|\mathbf{h}\|_{\ell_2}^2 + \lambda \left(\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \right)^q - \frac{\nu\beta}{\sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} + \frac{\nu^2}{2} - e_{q, \Sigma} \left(\frac{\beta}{\sqrt{n}} \Sigma^{-1/2} \mathbf{h} - \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\boldsymbol{\mu}}; \lambda \right) \right\} \\
& \quad + \frac{\alpha \tau_h}{2} \\
& = \min_{\tau_h, \lambda \geq 0, \nu} \frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2n} \|\mathbf{h}\|_{\ell_2}^2 + \lambda_0 \left(\frac{\tau_h \gamma_0}{\alpha} \right)^q - \frac{\nu\beta}{\sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} + \frac{\nu^2}{2} - e_{q, \Sigma} \left(\frac{\beta}{\sqrt{n}} \Sigma^{-1/2} \mathbf{h} - \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_p}^q \right) \right\} \\
& \quad (B.76) \\
& \quad + \frac{\alpha \tau_h}{2}
\end{aligned}$$

where we used the reparameterization $\lambda_0 := \frac{\lambda}{\|\boldsymbol{\mu}\|_{\ell_p}^q}$. Next we use Assumption 7 to take the limit of the above expression as $n \rightarrow \infty$. By definition of function E we have

$$\lim_{n \rightarrow \infty} e_{q, \Sigma} \left(\frac{\beta}{\sqrt{n}} \Sigma^{-1/2} \mathbf{h} - \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_p}^q \right) = \mathbb{E} \left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right); \lambda_0 \right).$$

Also, since $\mathbf{h} \sim \mathbf{N}(0, \mathbf{I}_d)$ we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} \|\mathbf{h}\|_{\ell_2}^2 = \frac{1}{\delta}, \quad \lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} = 0.$$

Using the above two equations in (B.76) we have

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \min_{\mathbf{w}} \left\{ \alpha \left\| \frac{\beta}{\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{h} + \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \Sigma^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{w} \right\|_{\ell_2} - \mathbf{w}^T \tilde{\boldsymbol{\mu}} \theta + \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \|\mathbf{w}\|_{\ell_p} \right\} \\
& \quad (B.77) \\
& = \frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^q + \frac{\nu^2}{2} - \mathbb{E} \left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right); \lambda_0 \right) \right\} + \frac{\alpha \tau_h}{2}.
\end{aligned}$$

Finally, incorporating the above equation in (B.62) and using Assumption 2, the AO problem simplifies to:

$$\begin{aligned} & \max_{0 \leq \beta \leq K} \min_{0 \leq \alpha, |\theta| \leq K', 0 < \gamma_0 < K'', 0 < \tau_g} \frac{\beta \tau_g}{2} + L \left(\sqrt{\alpha^2 + a^2 \theta^2}, V\theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta} \right) \\ & \quad (B.78) \\ & \quad - \min_{\tau_h, \lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^q + \frac{\nu^2}{2} - \mathbb{E} \left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right); \lambda_0 \right) \right\} + \frac{\alpha \tau_h}{2} \right] \end{aligned}$$

Now recall the argument after (B.73) that the objective is jointly convex in $(\alpha, \gamma, \theta, \tau_g)$ and jointly concave in (β, τ_h) . We used Lemma B.5 to provide alternative characterization for quantity $\|\mathbf{P}_B(\mathbf{h}) - \mathbf{h}\|_{\ell_2}^2$, which led into introducing the new variables λ_0, ν . Therefore, the objective (B.78), after maximization over λ_0, ν , is jointly convex in $(\alpha, \gamma_0, \theta, \tau_g)$ and jointly concave in (β, τ_h) . Because of that we can interchange the order of minimization and minimization over using Sion's minimax theorem to get the following.

$$\begin{aligned} & \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} \max_{0 \leq \beta, \tau_h} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) \\ & D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) = \frac{\beta \tau_g}{2} + L \left(\sqrt{\alpha^2 + a^2 \theta^2}, V\theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta} \right) \\ & \quad (B.79) \\ & \quad - \min_{\lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^q + \frac{\nu^2}{2} - \mathbb{E} \left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right); \lambda_0 \right) \right\} + \frac{\alpha \tau_h}{2} \right]. \end{aligned}$$

B.3.3. Uniform convergence of the auxiliary problem to its scalarized version

We showed that the auxiliary optimization objective converges pointwise to the function D_{ns} given by (B.79). However, we are interested in the minimax optimal solution of the auxiliary problem and need to have convergence of optimal points to the minimax solution of D . What is required for this aim is (local) uniform convergence of the auxiliary objective to function D . This can be shown by following similar arguments as in [63, Lemma A.5] that is essentially based on a result known as “convexity lemma” in the literature (see e.g. [39, Lemma 7.75]) by which pointwise convergence of convex functions implies uniform convergence in compact subsets.

B.3.4. Uniqueness of the solution of the AO problem

First note that since the loss $\ell(t)$ is a convex function and $\frac{1}{2\mu}(x-t)^2$ is jointly convex in (x, t, μ) , then $\frac{1}{2\mu}(x-t)^2 + \ell(t)$ is jointly convex in (x, t, μ) . Given that partial minimization preserves convexity, the Moreau envelope $e_\ell(x; \mu)$ is jointly convex in (x, μ) . In addition, by using the result of [63, Lemma 4.4] the

expected Moreau envelope of a convex function is jointly “strictly” convex (indeed this holds without requiring any strong or strict convexity assumption on the function itself). An application of this result to our case implies that $L(a, b, \mu)$ is jointly strictly convex in $\mathbb{R}_{\geq 0} \times \mathbb{R} \times \mathbb{R}_{\geq 0}$.

In addition, as we argued before the function D_{ns} given by (B.79) is jointly convex in $(\alpha, \gamma_0, \theta, \tau_g)$ and jointly concave in (β, τ_h) . Hence, using strict convexity of $L(a, b, \mu)$, the function D_{ns} is indeed jointly “strictly” convex in $(\alpha, \gamma_0, \theta, \tau_g)$ and jointly concave in (β, τ_h) .

As the next step, we note that $\max_{\beta, \tau_h} D(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h)$ is strictly convex in $(\alpha, \gamma_0, \theta, \tau_g)$. This follows from the fact that if a function $f(\mathbf{x}, \mathbf{y})$ is strictly convex in \mathbf{x} , then $\max_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is also strictly convex in \mathbf{x} . Moreover, by using the result of [63, Lemma C.5] we have that $\inf_{\tau_g > 0} \max_{\beta, \tau_h} D(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h)$ is strictly convex in $(\alpha, \gamma_0, \theta)$ and therefore has a unique minimizer $(\alpha_*, \gamma_{0*}, \theta_*)$. This concludes the part (a) of the theorem and the given scalar minimax optimization to characterize the limiting behavior of parameter of interest α, γ_0, θ .

Part (b) of the theorem follows readily from our definition of parameters α, θ and γ . Part (c) of the theorem also follows from combining Lemma 2.1 with part (b) of the theorem.

This completes the proof of Theorem 6.4.

B.4. Proof of Remark 6.1 We start by establishing an explicit expression for the weighted Moreau envelope $e_{q, \Sigma}$ for case of $p = q = 2$.

Lemma B.6 *We have*

$$e_{2, \Sigma}(\mathbf{x}; \lambda) = \lambda \left\| (\Sigma + 2\lambda \mathbf{I})^{-1/2} \Sigma^{1/2} \mathbf{x} \right\|_{\ell_2}^2$$

The proof of Lemma B.6 is given in Appendix E.9.

Suppose that items (i), (ii) in the statement of the remark are satisfied. We then prove that Assumption 6 and 7 hold.

Proof [Verification of Assumption 6] To check Assumption 6 for $p = q = 2$, we use Lemma B.6 to get

$$\begin{aligned} & \lim_{n \rightarrow \infty} e_{2, \mathbf{I} + b_0 \Sigma} \left((\mathbf{I} + b_0 \Sigma)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \tilde{\mu} \right\}; b_1 \|\mu\|_{\ell_2}^2 \right) \\ & \quad (\text{B.80}) \\ & = \lim_{n \rightarrow \infty} b_1 \|\mu\|_{\ell_2}^2 \left\| \left((1 + 2b_1 \|\mu\|_{\ell_2}^2) \mathbf{I} + b_0 \Sigma \right)^{-1/2} (\mathbf{I} + b_0 \Sigma)^{-1/2} \left(\frac{c_0}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \tilde{\mu} \right) \right\|_{\ell_2}^2 \end{aligned}$$

Consider a singular value decomposition $\Sigma = \mathbf{U} \mathbf{S} \mathbf{U}^T$ with $\mathbf{S} = \text{diag}(s_1, \dots, s_d)$, and the first column of \mathbf{U} being $\tilde{\mu}$ and $s_1 = a^2$ (Recall that $\tilde{\mu}$ is a singular

value of Σ with eigenvalue a^2 .) Also let $\tilde{\mathbf{h}} := \mathbf{U}^T \mathbf{h} \sim \mathbf{N}(0, \mathbf{I}_d)$. Continuing from (B.80) we write

$$\begin{aligned} & \lim_{n \rightarrow \infty} e_{2, \mathbf{I} + b_0 \Sigma} \left((\mathbf{I} + b_0 \Sigma)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \tilde{\mu} \right\}; b_1 \|\mu\|_{\ell_2}^2 \right) \\ & \quad (\text{B.81}) \\ & = \lim_{n \rightarrow \infty} \frac{1}{n} b_1 \|\mu\|_{\ell_2}^2 \left\| \mathbf{U} \left((1 + 2b_1 \|\mu\|_{\ell_2}^2) \mathbf{I} + b_0 \mathbf{S} \right)^{-1/2} (\mathbf{I} + b_0 \mathbf{S})^{-1/2} \mathbf{U}^T \left(\frac{c_0}{2} \Sigma^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \sqrt{n} \tilde{\mu} \right) \right\|_{\ell_2}^2 \end{aligned}$$

Write $\mathbf{U} = [\tilde{\mu}, \tilde{\mathbf{U}}]$ and $\tilde{\mathbf{S}} = \text{diag}(s_2, \dots, s_d)$. In addition, define $\tilde{\mathbf{h}} := \tilde{\mathbf{U}}^T \mathbf{h} \sim \mathbf{N}(0, \mathbf{I}_{d-1})$. We then have

$$\mathbf{U}^T \Sigma^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} = \begin{pmatrix} 0 \\ \tilde{\mathbf{S}}^{1/2} \tilde{\mathbf{h}} \end{pmatrix}, \quad \mathbf{U}^T \tilde{\mu} = \mathbf{e}_1, \quad \lim_{n \rightarrow \infty} \|\mu\|_{\ell_2} = \sigma_{M,2},$$

in probability, with the last limit following from item (i) in the statement Remark 6.1. Using the above identities in (B.81) we get

$$\begin{aligned} & \lim_{n \rightarrow \infty} e_{2, \mathbf{I} + b_0 \Sigma} \left((\mathbf{I} + b_0 \Sigma)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \Sigma^{1/2} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \tilde{\mu} \right\}; b_1 \|\mu\|_{\ell_2}^2 \right) \\ & = \lim_{n \rightarrow \infty} \frac{b_1 \sigma_{M,2}^2}{n} \left\{ \frac{c_1^2 n}{4} \cdot \frac{1}{(1 + 2b_1 \sigma_{M,2}^2 + b_0 a^2)(1 + b_0 a^2)} + \frac{c_0^2}{4} \sum_{i=2}^d \frac{s_i \tilde{h}_i^2}{(1 + b_0 s_i)(1 + 2b_1 \sigma_{M,2}^2 + b_0 s_i)} \right\} \\ & \quad (\text{B.82}) \\ & = b_1 \sigma_{M,2}^2 \left\{ \frac{c_1^2}{4(1 + b_0 a^2)(1 + 2b_1 \sigma_{M,2}^2 + b_0 a^2)} + \frac{c_0^2}{4\delta} \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=2}^d \frac{s_i \tilde{h}_i^2}{(1 + b_0 s_i)(1 + 2b_1 \sigma_{M,2}^2 + b_0 s_i)} \right\} \end{aligned}$$

Define $\nu_i := s_i(1 + b_0 s_i)^{-1}(1 + 2b_1 \sigma_{M,2}^2 + b_0 s_i)^{-1}$. Then the last sum reads as $\frac{1}{d} \sum_{i=2}^d \nu_i \tilde{h}_i^2$. Recall that $\tilde{\mathbf{h}} \sim \mathbf{N}(0, \mathbf{I}_{d-1})$. Therefore, by applying the Kolmogorov's criterion of SLLN the above limit exists (almost surely and so in probability as well) provided that $\frac{1}{d^2} \sum_{i=2}^d \nu_i^2 \text{Var}(\tilde{h}_i^2) < \infty$. We note that $\text{Var}(\tilde{h}_i^2) = 2$ and since $\nu_i \geq 0$, we have

$$\frac{1}{d^2} \sum_{i=2}^d \nu_i^2 \leq \left(\frac{1}{d} \sum_{i=2}^d \nu_i \right)^2.$$

Hence it suffices to show that $\frac{1}{d} \sum_{i=2}^d \nu_i < \infty$. Now by item (ii) of Remark 6.1, the empirical distribution of eigenvalues of Σ converges weakly to a distribu-

tion ρ with Stieltjes transform $S_\rho(z) := \int \frac{\rho(t)}{z-t} dt$. We write

$$\begin{aligned}
 & \frac{s_i}{(1+b_0 s_i)(1+2b_1 \sigma_{M,2}^2 + b_0 s_i)} \\
 &= \frac{1}{2b_0 b_1 \sigma_{M,2}^2} \left\{ -\frac{1}{1+b_0 s_i} + \frac{1+2b_1 \sigma_{M,2}^2}{1+2b_1 \sigma_{M,2}^2 + b_0 s_i} \right\} \\
 (B.83) \quad &= \frac{1}{2b_0^2 b_1 \sigma_{M,2}^2} \left\{ -\frac{1}{\frac{1}{b_0} + s_i} + \frac{1+2b_1 \sigma_{M,2}^2}{\frac{1+2b_1 \sigma_{M,2}^2}{b_0} + s_i} \right\}
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=2}^d \nu_i &= \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=2}^d \frac{s_i}{(1+b_0 s_i)(1+2b_1 \sigma_{M,2}^2 + b_0 s_i)} \\
 &= \frac{1}{2b_0^2 b_1 \sigma_{M,2}^2} \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=2}^d \left\{ -\frac{1}{\frac{1}{b_0} + s_i} + \frac{1+2b_1 \sigma_{M,2}^2}{\frac{1+2b_1 \sigma_{M,2}^2}{b_0} + s_i} \right\} \\
 (B.84) \quad &= \frac{1}{2b_0^2 b_1 \sigma_{M,2}^2} \left\{ S_\rho \left(-\frac{1}{b_0} \right) - (1+2b_1 \sigma_{M,2}^2) S_\rho \left(-\frac{1+2b_1 \sigma_{M,2}^2}{b_0} \right) \right\}.
 \end{aligned}$$

It is worth noting that although the sum is over $2 \leq i \leq d$, the term for $i = 1$ is $O(1/d)$ and is negligible in the limit. Therefore, we can include that in our calculation above. By using Equation (B.84) in (B.82) we get that Assumption 6 holds with

$$\begin{aligned}
 F(c_0, c_1; b_0, b_1) &= \frac{b_1 \sigma_{M,2}^2 c_1^2}{4(1+b_0 a^2)(1+2b_1 \sigma_{M,2}^2 + b_0 a^2)} \\
 (B.85) \quad &+ \frac{b_1 \sigma_{M,2}^2 c_0^2}{8\delta b_0^2 b_1 \sigma_{M,2}^2} \left\{ S_\rho \left(-\frac{1}{b_0} \right) - (1+2b_1 \sigma_{M,2}^2) S_\rho \left(-\frac{1+2b_1 \sigma_{M,2}^2}{b_0} \right) \right\}.
 \end{aligned}$$

■

Proof [Verification of Assumption 7] To check Assumption 7 we use Lemma B.6 and write

$$\begin{aligned}
 & \lim_{n \rightarrow \infty} e_{2,\Sigma} \left(\frac{c_0}{\sqrt{n}} \Sigma^{-1/2} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \right) \\
 &= \lim_{n \rightarrow \infty} \lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \left\| (\Sigma + 2\lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \mathbf{I})^{-1/2} \Sigma^{1/2} \left(\frac{c_0}{\sqrt{n}} \Sigma^{-1/2} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}} \right) \right\|_{\ell_2}^2 \\
 (B.86) \quad &= \lim_{n \rightarrow \infty} \lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \left\| (\Sigma + 2\lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \mathbf{I})^{-1/2} \left(\frac{c_0}{\sqrt{n}} \mathbf{h} - c_1 a \tilde{\boldsymbol{\mu}} \right) \right\|_{\ell_2}^2
 \end{aligned}$$

Consider a singular value decomposition $\mathbf{\Sigma} = \mathbf{U} \mathbf{S} \mathbf{U}^T$ with $\mathbf{S} = \text{diag}(s_1, \dots, s_d)$, and the first column of \mathbf{U} being $\tilde{\boldsymbol{\mu}}$ and $s_1 = a^2$ (Recall that $\tilde{\boldsymbol{\mu}}$ is a singular value of $\mathbf{\Sigma}$ with eigenvalue a^2 .) Also let $\tilde{\mathbf{h}} := \mathbf{U}^T \mathbf{h} \sim \mathcal{N}(0, \mathbf{I}_d)$. Continuing from (B.86) we write

$$\begin{aligned}
 & \lim_{n \rightarrow \infty} e_{2, \mathbf{\Sigma}} \left(\frac{c_0}{\sqrt{n}} \mathbf{\Sigma}^{-1/2} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \right) \\
 &= \lim_{n \rightarrow \infty} \frac{1}{n} \lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \left\| \mathbf{U} (\mathbf{S} + 2\lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \mathbf{I})^{-1/2} (c_0 \tilde{\mathbf{h}} - c_1 \sqrt{n} a \mathbf{e}_1) \right\|_{\ell_2}^2 \\
 &= \lim_{n \rightarrow \infty} \frac{\lambda_0 \sigma_{M,2}^2}{n} \left\{ \frac{(c_0 \tilde{h}_1 - c_1 \sqrt{n} a)^2}{a^2 + 2\lambda_0 \sigma_{M,2}^2} + \sum_{i=2}^d \frac{c_0^2 \tilde{h}_i^2}{s_i + 2\lambda_0 \sigma_{M,2}^2} \right\} \\
 \text{(B.87)} \quad &= \lambda_0 \sigma_{M,2}^2 \left\{ \frac{c_1^2 a^2}{a^2 + 2\lambda_0 \sigma_{M,2}^2} + \frac{1}{\delta} \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=2}^d \frac{c_0^2 \tilde{h}_i^2}{s_i + 2\lambda_0 \sigma_{M,2}^2} \right\}.
 \end{aligned}$$

By applying the Kolmogorov's criterion of SLLN the above limit exists (almost surely and so in probability as well) provided that $\frac{1}{d^2} \sum_{i=2}^d \frac{1}{(s_i + 2\lambda_0 \sigma_{M,2}^2)^2} < \infty$. Note that since $\lambda_0, s_i \geq 0$, we have

$$\frac{1}{d^2} \sum_{i=2}^d \frac{1}{(s_i + 2\lambda_0 \sigma_{M,2}^2)^2} \leq \frac{1}{d^2} \sum_{i=2}^d \frac{1}{4\lambda_0^2 \sigma_{M,2}^4} \rightarrow 0.$$

By using the LLN we obtain that the summation in (B.87) converges (almost surely) to its expectation. Now recalling item (ii) in Remark 6.1, we know that the empirical distribution of eigenvalues of $\mathbf{\Sigma}$ converges weakly to a distribution ρ with Stieltjes transform $S_\rho(z) := \int \frac{\rho(t)}{z-t} dt$, and therefore we have

$$\begin{aligned}
 \text{E}(c_0, c_1; \lambda_0) &:= \lim_{n \rightarrow \infty} e_{2, \mathbf{\Sigma}} \left(\frac{c_0}{\sqrt{n}} \mathbf{\Sigma}^{-1/2} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_2}^2 \right) \\
 \text{(B.88)} \quad &= \lambda_0 \sigma_{M,2}^2 \left\{ \frac{c_1^2 a^2}{a^2 + 2\lambda_0 \sigma_{M,2}^2} - \frac{c_0^2}{\delta} S_\rho(-2\lambda_0 \sigma_{M,2}^2) \right\}.
 \end{aligned}$$

■

APPENDIX C: PROOFS FOR ISOTROPIC GAUSSIAN MODEL (SECTION 4)

This section is devoted to the proof of our theorems for the isotropic Gaussian model. We discuss how these theorems can be derived as special cases of our results for the anisotropic model, after some algebraic simplifications.

The claim of Theorem 4.1 on the separability threshold is an immediate corollary of Theorem 6.1, with $\Sigma = \mathbf{I}_{p \times p}$ and $a = 1$. We next move to the two other theorems on precise characterization of standard and robust accuracy in the separable and non-separable regimes.

C.1. Proof of Theorem 4.3 Suppose that Assumption 4 in the statement of Theorem 4.3 holds. We first show that this assumption implies Assumption 6, required by Theorem 6.3, in case of $\Sigma = \mathbf{I}$ and then show how Theorem 4.3 can be derived as a special case of Theorem 6.3.

To prove Assumption 6(b) for isotropic case, we use the following two properties of the weighted Moreau envelop that holds for all $q \geq 0$:

$$(C.1) \quad e_{q, \alpha \mathbf{I}}(\mathbf{x}; \lambda) = \alpha e_{q, \mathbf{I}}\left(\mathbf{x}, \frac{\lambda}{\alpha}\right),$$

$$(C.2) \quad \frac{1}{b^2} e_{q, \mathbf{I}}\left(b\mathbf{x}; \frac{\lambda}{b^{q-2}}\right) = e_{q, \mathbf{I}}(\mathbf{x}; \lambda).$$

Combining the above two identities we get

$$e_{q, \alpha \mathbf{I}}(\alpha^{-1} \mathbf{x}; \lambda) = \alpha e_{q, \mathbf{I}}\left(\frac{\mathbf{x}}{\alpha}; \frac{\lambda}{\alpha}\right) = \frac{\alpha}{b^2} e_{q, \mathbf{I}}\left(b \frac{\mathbf{x}}{\alpha}; \frac{\lambda}{\alpha b^{q-2}}\right).$$

Using the above identity with $\alpha = 1 + b_0$ and $b = \alpha \sqrt{d}$ we have

$$(C.3) \quad \begin{aligned} & e_{q, (1+b_0) \mathbf{I}}\left((1+b_0)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \tilde{\mu} \right\}; b_1 \|\mu\|_{\ell_p}^q\right) \\ &= \frac{1}{(1+b_0)d} e_{q, \mathbf{I}}\left(\frac{\sqrt{d}c_0}{2\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1 \sqrt{d}}{2} \tilde{\mu}; \frac{b_1 \|\mu\|_{\ell_p}^q}{(1+b_0)^{q-1} d^{\frac{q}{2}-1}}\right) \end{aligned}$$

We next proceed to take the limit of the above expression as $n \rightarrow \infty$. By Assumption 4 we have

$$(C.4) \quad \|\mu\|_{\ell_p}^q \rightarrow \sigma_{M,p}^q d^{\frac{q}{2}-1}, \quad \|\mu\|_{\ell_2} \rightarrow \sigma_{M,2},$$

with high probability. Also by Assumption 1, we have $n/d \rightarrow \delta$. Therefore,

$$(C.5) \quad \begin{aligned} & \lim_{n \rightarrow \infty} e_{q, (1+b_0) \mathbf{I}}\left((1+b_0)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1}{2} \tilde{\mu} \right\}; b_1 \|\mu\|_{\ell_p}^q\right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{(1+b_0)d} e_{q, \mathbf{I}}\left(\frac{\sqrt{d}c_0}{2\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1 \sqrt{d}}{2} \tilde{\mu}; \frac{b_1 \|\mu\|_{\ell_p}^q}{(1+b_0)^{q-1} d^{\frac{q}{2}-1}}\right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{(1+b_0)d} e_{q, \mathbf{I}}\left(\frac{\sqrt{d}c_0}{2\sqrt{n}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1 \sqrt{d} \mu}{2 \|\mu\|_{\ell_2}}; \frac{b_1 \|\mu\|_{\ell_p}^q}{(1+b_0)^{q-1} d^{\frac{q}{2}-1}}\right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{(1+b_0)d} e_{q, \mathbf{I}}\left(\frac{c_0}{2\sqrt{\delta}} \mathbf{P}_{\mu}^{\perp} \mathbf{h} - \frac{c_1 \sqrt{d} \mu}{2 \sigma_{M,2}}; b_1 (1+b_0)^{1-q} \sigma_{M,p}^q\right) \end{aligned}$$

We next note that by definition of the weighted Moreau envelop we have $e_{q,\mathbf{I}}(\mathbf{x}; \lambda) = \sum_{i=1}^d J_q(x_i; \lambda)$. Also, by Assumption 7 the empirical distribution of entries of $\sqrt{d}\boldsymbol{\mu}$ converges weakly to distribution \mathbb{P}_M . Therefore, continuing from (C.5) we can write

$$\begin{aligned}
& \lim_{n \rightarrow \infty} e_{q,(1+b_0)\mathbf{I}} \left((1+b_0)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \mathbf{P}_{\boldsymbol{\mu}}^\perp \mathbf{h} - \frac{c_1}{2} \tilde{\boldsymbol{\mu}} \right\}; b_1 \|\boldsymbol{\mu}\|_{\ell_p}^q \right) \\
& \stackrel{(a)}{=} \lim_{n \rightarrow \infty} e_{q,(1+b_0)\mathbf{I}} \left((1+b_0)^{-1} \left\{ \frac{c_0}{2\sqrt{n}} \mathbf{h} - \frac{c_1}{2} \tilde{\boldsymbol{\mu}} \right\}; b_1 \|\boldsymbol{\mu}\|_{\ell_p}^q \right) \\
& = \lim_{n \rightarrow \infty} \frac{1}{(1+b_0)d} \sum_{i=1}^d J_q \left(\frac{c_0 h_i}{2\sqrt{\delta}} - \frac{c_1}{2} \frac{\sqrt{d}\mu_i}{\sigma_{M,2}}; b_1 (1+b_0)^{1-q} \sigma_{M,p}^q \right) \\
& \stackrel{(b)}{=} \frac{1}{1+b_0} \mathbb{E} \left[J_q \left(\frac{c_0 h}{2\sqrt{\delta}} - \frac{c_1 M}{2\sigma_{M,2}}; b_1 (1+b_0)^{1-q} \sigma_{M,p}^q \right) \right] = (1+b_0)^{-1} \mathcal{J} \left(\frac{c_0}{2}, \frac{c_1}{2}; b_1 (1+b_0)^{1-q} \right),
\end{aligned} \tag{C.6}$$

where the expectation in (b) is taken with respect to the independent random variables $h \sim \mathbf{N}(0, 1)$ and $M \sim \mathbb{P}_M$. The last equality follows by definition of function \mathcal{J} given by (4.7) and by deploying Assumption 4.

Here, (a) follows by writing

$$\frac{c_0}{2\sqrt{\delta}} \mathbf{P}_{\boldsymbol{\mu}}^\perp \mathbf{h} - \frac{c_1 \sqrt{d} \boldsymbol{\mu}}{2\sigma_{M,2}} = \frac{c_0}{2\sqrt{\delta}} \mathbf{h} - \left(\frac{c_1 \sqrt{d}}{2\sigma_{M,2}} + \frac{c_0}{2\sqrt{\delta} \|\boldsymbol{\mu}\|_{\ell_2}} \mathbf{h}^T \tilde{\boldsymbol{\mu}} \right) \boldsymbol{\mu}$$

and noting that $\tilde{\boldsymbol{\mu}}^T \mathbf{h} \sim \mathbf{N}(0, 1)$ since $\|\tilde{\boldsymbol{\mu}}\|_{\ell_2} = 1$, and $\|\boldsymbol{\mu}\|_{\ell_2} \rightarrow \sigma_{M,2}$ which implies that the last term in the right-hand side is dominated by the second term therein that is of order \sqrt{d} .

The chain of equalities in (C.6) shows that Assumption 6(a) is satisfied by $F(c_0, c_1; b_0, b_1) = \mathcal{J} \left(\frac{c_0}{2}, \frac{c_1}{2}; b_1 (1+b_0)^{1-q} \right)$ for the isotropic model.

Assumption 6(b) also clearly holds for isotropic model ($\boldsymbol{\Sigma} = \mathbf{I}$) with $S_\rho(z) = \frac{1}{z-1}$.

Now that Assumption 6 holds we can use the result of Theorem 6.3 for the special case of $\boldsymbol{\Sigma} = \mathbf{I}$. As we showed above for this case, we have the following identities

$$\tag{C.7} \quad F(c_0, c_1; b_0, b_1) = \mathcal{J} \left(\frac{c_0}{2}, \frac{c_1}{2}; b_1 (1+b_0)^{1-q} \right), \quad S_\rho(z) = \frac{1}{z-1}.$$

Now by using these identities in the AO problem (6.2) and after some simple algebraic manipulation we obtain the AO problem (4.8).

C.2. Proof of Theorem 4.5 We prove Theorem 4.5 as a special case of Theorem 6.4. We first show that in the isotropic case, Assumption 4 implies Assumption 7, required by Theorem 6.4.

Note that by Assumption 4 we have

$$(C.8) \quad \|\boldsymbol{\mu}\|_{\ell_p}^q \rightarrow \sigma_{M,p}^q d^{\frac{q}{2}-1}, \quad \|\boldsymbol{\mu}\|_{\ell_2} \rightarrow \sigma_{M,2},$$

with high probability.

We then write

$$\begin{aligned} \lim_{n \rightarrow \infty} e_{q,\mathbf{I}} \left(\frac{c_0}{\sqrt{n}} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_p}^q \right) &\stackrel{(a)}{=} \lim_{n \rightarrow \infty} \frac{1}{d} e_{q,\mathbf{I}} \left(\frac{c_0 \sqrt{d}}{\sqrt{n}} \mathbf{h} - c_1 \sqrt{d} \tilde{\boldsymbol{\mu}}; \frac{\lambda_0 \|\boldsymbol{\mu}\|_{\ell_p}^q}{d^{\frac{q}{2}-1}} \right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{d} e_{q,\mathbf{I}} \left(\frac{c_0 \sqrt{d}}{\sqrt{n}} \mathbf{h} - c_1 \frac{\sqrt{d} \boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_{\ell_2}}; \frac{\lambda_0 \|\boldsymbol{\mu}\|_{\ell_p}^q}{d^{\frac{q}{2}-1}} \right) \\ &\stackrel{(b)}{=} \lim_{n \rightarrow \infty} \frac{1}{d} e_{q,\mathbf{I}} \left(\frac{c_0}{\sqrt{\delta}} \mathbf{h} - c_1 \frac{\sqrt{d} \boldsymbol{\mu}}{\sigma_{M,2}}; \lambda_0 \sigma_{M,p}^q \right) \\ &\stackrel{(c)}{=} \lim_{n \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d J_q \left(\frac{c_0 h_i}{\sqrt{\delta}} - c_1 \frac{\sqrt{d} \mu_i}{\sigma_{M,2}}; \lambda_0 \sigma_{M,p}^q \right) \\ &\stackrel{(d)}{=} \mathbb{E} \left[J_q \left(\frac{c_0}{\sqrt{\delta}} h - c_1 \frac{M}{\sigma_{M,2}}; \lambda_0 \sigma_{M,p}^q \right) \right] \\ (C.9) \quad &= \mathcal{J}(c_0, c_1; \lambda_0). \end{aligned}$$

Here (a) follows from (C.2) with $b = \sqrt{d}$; (b) follows from (C.8); (c) holds due to the identity $e_{q,\mathbf{I}}(\mathbf{x}; \lambda) = \sum_{i=1}^d J_q(x_i; \lambda)$, which follows readily from the definition of weighted Moreau envelop $e_{q,\mathbf{I}}$ and the function J_q given by (4.5). Finally, (d) holds due to Assumption 4. The series of equalities (C.9) implies that Assumption 7 holds in isotropic case with

$$(C.10) \quad \mathbb{E}(c_0, c_1; \lambda_0) = \mathcal{J}(c_0, c_1; \lambda_0).$$

Having Assumption 6 in place, we can specialize the result of Theorem 6.4 to isotropic model. Substituting for $\mathbb{E}(c_0, c_1; \lambda_0)$ from (C.10) in the AO problem (6.4) yields the AO problem (4.16).

APPENDIX D: PROOFS FOR SPECIAL CASES OF P (SECTION 5)

D.1. Proof of Corollary 5.1 Part (a) is already proved in Example 2, cf. (4.4).

Proof [Part (b)] We start by an explicit characterization of \mathcal{J} function for case of $p = q = 2$.

Lemma D.1 Recall function $\mathcal{J}(c_0, c_1; \lambda_0)$ given by

$$(D.1) \quad \mathcal{J}(c_0, c_1; \lambda_0) = \mathbb{E} \left[J_q \left(\frac{c_0}{\sqrt{\delta}} h - c_1 \frac{M}{\sigma_{M,2}}; \lambda_0 \sigma_{M,p}^q \right) \right],$$

Then the following identity holds for case of $p = q = 2$:

$$\mathcal{J}(c_0, c_1; \lambda_0) = \frac{\lambda}{\alpha^2 + 2\lambda\alpha} \|\mathbf{x}\|_{\ell_2}^2$$

Proof It is straightforward to see that

$$J_2(x; \lambda) = \frac{\lambda}{1 + 2\lambda} x^2.$$

Therefore,

$$\mathcal{J}(c_0, c_1; \lambda_0) = \frac{\lambda_0 \sigma_{M,2}^2}{1 + 2\lambda_0 \sigma_{M,2}^2} \mathbb{E} \left[\left(\frac{c_0}{\sqrt{\delta}} h - c_1 \frac{M}{\sigma_{M,2}} \right)^2 \right] = \frac{\lambda_0 \sigma_{M,2}^2}{1 + 2\lambda_0 \sigma_{M,2}^2} \left(\frac{c_0^2}{\delta} + c_1^2 \right).$$

■

Using Lemma D.1 in AO problem (4.8) for $q = 2$, we have

$$(D.2) \quad \begin{aligned} D_s(\alpha, \gamma_0, \theta, \beta, \lambda_0, \eta, \tilde{\eta}) &= 2 \left(1 + \frac{\eta}{2\alpha} \right)^{-1} \mathcal{J} \left(\frac{\beta}{2}, \frac{\tilde{\eta}}{2}; \frac{\lambda_0}{2\gamma_0} \left(1 + \frac{\eta}{2\alpha} \right)^{-1} \right) \\ &\quad - \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right) \frac{1}{4(1 + \frac{\eta}{2\alpha})} - \lambda_0 \gamma_0 - \frac{\eta\alpha}{2} - \tilde{\eta}\theta \\ &\quad + \beta \sqrt{\mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right]} \\ &= \frac{\frac{\lambda_0}{4\gamma_0} \sigma_{M,2}^2}{(1 + \frac{\eta}{2\alpha})^2 + (1 + \frac{\eta}{2\alpha}) \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2} \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right) \\ &\quad - \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right) \frac{1}{4(1 + \frac{\eta}{2\alpha})} - \lambda_0 \gamma_0 - \frac{\eta\alpha}{2} - \tilde{\eta}\theta \\ &\quad + \beta \sqrt{\mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right]} \\ &= - \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right) \frac{1}{4(1 + \frac{\eta}{2\alpha} + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2)} - \lambda_0 \gamma_0 - \frac{\eta\alpha}{2} - \tilde{\eta}\theta \\ &\quad + \beta \sqrt{\mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right]}. \end{aligned}$$

Setting $\frac{\partial D_s}{\partial \beta}$ to zero we conclude that

$$\widehat{\beta} = 2\delta \left(1 + \frac{\eta}{2\alpha} + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2\right) \sqrt{\mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right]}.$$

Thus the AO problem reduces to

$$(D.3) \quad \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\lambda_0, \eta \geq 0, \tilde{\eta}} -\tilde{\eta}^2 \frac{1}{4(1 + \frac{\eta}{2\alpha} + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2)} - \lambda_0 \gamma_0 - \frac{\eta \alpha}{2} - \tilde{\eta} \theta + \delta \left(1 + \frac{\eta}{2\alpha} + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2\right) \mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right].$$

Setting the derivative with respect to $\tilde{\eta}$ to zero we arrive at

$$\widehat{\tilde{\eta}} = -2\theta \left(1 + \frac{\eta}{2\alpha} + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2\right),$$

which further simplifies the AO problem to

$$(D.4) \quad \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\lambda_0, \eta \geq 0} -\lambda_0 \gamma_0 - \frac{\eta \alpha}{2} + \theta^2 \left(1 + \frac{\eta}{2\alpha} + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2\right) + \delta \left(1 + \frac{\eta}{2\alpha} + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2\right) \mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right].$$

Note that if $\alpha^2 < \theta^2 + \delta \mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right]$ then the maximum over η is $+\infty$. Furthermore, when $\alpha^2 \geq \theta^2 + \delta \mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right]$ then the optimal $\eta = 0$. Thus the above AO is equivalent to

$$(D.5) \quad \begin{aligned} & \min_{\alpha, \gamma_0 \geq 0, \theta} \max_{\lambda_0 \geq 0} -\lambda_0 \gamma_0 + \theta^2 \left(1 + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2\right) \\ & \quad + \delta \left(1 + \frac{\lambda_0}{\gamma_0} \sigma_{M,2}^2\right) \mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right] \\ & \text{subject to } \alpha^2 \geq \theta^2 + \delta \mathbb{E} \left[((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g)_+^2 \right] \end{aligned}$$

Using a similar argument for optimization over λ_0 , it is straightforward to see that the above optimization is equivalent to

$$\begin{aligned}
 (D.6) \quad & \min_{\alpha, \gamma_0 \geq 0, \theta} \theta^2 + \delta \mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right] \\
 & \text{subject to } \alpha^2 \geq \theta^2 + \delta \mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right], \\
 & \text{and } \frac{\gamma_0^2}{\sigma_{M,2}^2} \geq \theta^2 + \delta \mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right].
 \end{aligned}$$

Since the objective function is increasing in α and γ , then the optimal α and γ should make the inequality constraints equality and therefore $\gamma_0 = \alpha \sigma_{M,2}$. This brings us to the following problem:

$$\begin{aligned}
 (D.7) \quad & \min_{\alpha \geq 0, u} \alpha^2 \\
 & \text{subject to } \alpha^2 \geq \theta^2 + \delta \mathbb{E} \left[\left((1 + (\varepsilon_0 \alpha - \theta) \sigma_{M,2} + \alpha g) \right)_+^2 \right].
 \end{aligned}$$

By change of variable $u = \frac{\theta}{\alpha}$ we have

$$\begin{aligned}
 (D.8) \quad & \min_{\alpha \geq 0, u} \alpha^2 \\
 & \text{subject to } 1 \geq u^2 + \delta \mathbb{E} \left[\left(\frac{1}{\alpha} + (\varepsilon_0 - u) \sigma_{M,2} + g \right)_+^2 \right]
 \end{aligned}$$

By another change of variable $\tilde{\alpha} = \left(\frac{1}{\alpha} + \varepsilon_0 \sigma_{M,2} \right)^{-1}$ we have

$$\begin{aligned}
 (D.9) \quad & \min_{\frac{1}{\varepsilon_0 \sigma_{M,2}} \geq \tilde{\alpha} \geq 0, u} \left(\frac{1}{\tilde{\alpha}} - \varepsilon_0 \sigma_{M,2} \right)^{-2} \\
 & \text{subject to } 1 \geq u^2 + \delta \mathbb{E} \left[\left(\frac{1}{\tilde{\alpha}} - u \sigma_{M,2} + g \right)_+^2 \right]
 \end{aligned}$$

Since objective is increasing in $\tilde{\alpha}$ this is equivalent to

$$\begin{aligned}
 (D.10) \quad & \min_{\frac{1}{\varepsilon_0 \sigma_{M,2}} \geq \tilde{\alpha} \geq 0, u} \tilde{\alpha}^2 \\
 & \text{subject to } 1 \geq u^2 + \delta \mathbb{E} \left[\left(\frac{1}{\tilde{\alpha}} - u \sigma_{M,2} + g \right)_+^2 \right]
 \end{aligned}$$

Note that we can drop the constraint $\frac{1}{\varepsilon_0 \sigma_{M,2}} \geq \tilde{\alpha}$ because for $\frac{1}{\varepsilon_0 \sigma_{M,2}} = \tilde{\alpha}$ one can already find u that satisfies the inequality constraint. As such the optimal

$\tilde{\alpha}$ should be less than $\frac{1}{\varepsilon_0 \sigma_{M,2}}$. To see why, by letting $u = \frac{\theta}{\sqrt{1+\theta^2}}$ with θ the minimizer in separability condition (4.4) we have

$$\begin{aligned} u^2 + \delta \mathbb{E} \left[\left(\frac{1}{\tilde{\alpha}} - u \sigma_{M,2} + g \right)_+^2 \right] &= \frac{\theta^2}{1+\theta^2} + \delta \mathbb{E} \left[\left(\left(\varepsilon_0 - \frac{\theta}{\sqrt{1+\theta^2}} \right) \sigma_{M,2} + g \right)_+^2 \right] \\ &= \frac{\theta^2}{1+\theta^2} + \frac{\delta}{1+\theta^2} \mathbb{E} \left[\left((\sqrt{1+\theta^2} \varepsilon_0 - \theta) \sigma_{M,2} + \sqrt{1+\theta^2} g \right)_+^2 \right] \\ &\leq \frac{\theta^2}{1+\theta^2} + \frac{1}{1+\theta^2} = 1. \end{aligned}$$

This brings us to the following AO problem:

$$\begin{aligned} &\min_{\tilde{\alpha} \geq 0, u} \tilde{\alpha}^2 \\ \text{(D.11)} \quad &\text{subject to} \quad 1 \geq u^2 + \delta \mathbb{E} \left[\left(\frac{1}{\tilde{\alpha}} - u \sigma_{M,2} + g \right)_+^2 \right] \end{aligned}$$

Denoting by $\tilde{\alpha}_*$ the solution of the above problem, it is clear that by our change of variable we have

$$\text{(D.12)} \quad \alpha_* = (\tilde{\alpha}_*^{-1} - \varepsilon_0 \sigma_{M,2})^{-1}, \quad \theta_* = u_* \alpha_*, \quad \gamma_{0*} = \alpha_* \sigma_{M,2}.$$

This concludes the proof of part (b). \blacksquare

Proof [Part (c)] We focus on part of the AO problem (4.16) that involves the variables λ_0, ν, τ_h and specialize it to the case of $q = 2$:

$$\min_{\lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^2 + \frac{\nu^2}{2} - \mathcal{J} \left(\beta, \left(\frac{\tau_h \theta}{\alpha} + \nu \right); \lambda_0 \right) \right\} + \frac{\alpha \tau_h}{2} \right].$$

We next plug in for $\mathcal{J}(c_0, c_1; \lambda_0)$ from Lemma D.1 which results in

$$\min_{\lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} \right)^2 + \frac{\nu^2}{2} - \frac{\lambda_0 \sigma_{M,2}^2}{1 + 2\lambda_0 \sigma_{M,2}^2} \left(\frac{\beta^2}{\delta} + \left(\frac{\tau_h \theta}{\alpha} + \nu \right)^2 \right) \right\} + \frac{\alpha \tau_h}{2} \right].$$

Writing the first order optimality for λ_0, ν, τ_h we get a set of equations that admits a solution only if $\gamma_0 = \sigma_{M,2} \sqrt{\alpha^2 + \theta^2}$. Then,

$$\nu = 2\lambda_0 \sigma_{M,2}^2 \frac{\tau_h \theta}{\alpha}, \quad \tau_h = \frac{1}{1 + 2\lambda_0 \sigma_{M,2}^2} \frac{\beta}{\sqrt{\delta}}.$$

In this case, the value of λ_0 does not matter and the above part of the AO simplifies to $\alpha \beta / \sqrt{\delta}$.

This simplifies the AO problem (4.16) to

$$(D.13) \quad \max_{0 \leq \beta} \min_{\theta, 0 \leq \alpha, \tau_g} \frac{\beta \tau_g}{2} + L\left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2}\theta - \varepsilon_0 \sqrt{\alpha^2 + \theta^2}, \frac{\tau_g}{\beta}\right) - \frac{\alpha\beta}{\sqrt{\delta}}.$$

We next further simplifies the AO problem by solving for τ_g . We use the shorthand $L'_3(a, b; \mu) = \frac{\partial L}{\partial \mu} L(a, b; \mu)$ to denote the derivative of the expected Moreau envelop with respect to its third argument. Writing the first order optimality condition for β and τ_g in optimization (D.13), we get

$$(D.14) \quad \begin{aligned} \frac{\beta}{2} + \frac{1}{\beta} L'_3\left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2}\theta - \varepsilon_0 \sqrt{\alpha^2 + \theta^2}, \frac{\tau_g}{\beta}\right) &= 0, \\ \frac{\tau_g}{2} - \frac{\tau_g}{\beta^2} L'_3\left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2}\theta - \varepsilon_0 \sqrt{\alpha^2 + \theta^2}, \frac{\tau_g}{\beta}\right) - \frac{\alpha}{\sqrt{\delta}} &= 0. \end{aligned}$$

Combining the above two equations, we obtain $\tau_g = \frac{\alpha}{\sqrt{\delta}}$. Substituting for τ_g in (D.13), the AO problem for case of $q = 2$ simplifies to

$$(D.15) \quad \max_{0 \leq \beta} \min_{\theta, 0 \leq \alpha} L\left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2}\theta - \varepsilon_0 \sqrt{\alpha^2 + \theta^2}, \frac{\alpha}{\beta\sqrt{\delta}}\right) - \frac{\alpha\beta}{2\sqrt{\delta}}.$$

This completes the proof of part (c). ■

D.2. Proof of Theorem 5.2 Proof [Part (a)] The first part of the theorem is on precise characterization of the separability threshold. We use the result of Theorem 4.1 that holds for any choice of (p, q) , in particular $(p = 1, q = \infty)$, and relates the separability threshold to the spherical width. What is remaining to prove is the characterization of the spherical width given by (5.9). To this end, we follow a similar argument as in Lemma 4.2. However, since $\|\cdot\|_{\ell_q}^q$ is not well defined for $q = \infty$ (recall that ℓ_q is the dual norm of ℓ_p and $p = 1$), it requires a slightly different analysis. Specifically, in the Lagrangian we write the constraint $\|\mathbf{u}\|_{\ell_\infty} \leq \frac{1}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_\infty}}$ as the term $\|\mathbf{u}\|_{\ell_\infty} - \frac{1}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}}$, as compared to the case of finite q where we raised the both sides to power q to use the separability property of function $\|\cdot\|_{\ell_q}^q$. Then, by following a similar derivation as in (E.4), we obtain

$$(D.16) \quad \begin{aligned} & \min_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})} - \frac{1}{\sqrt{n}} \mathbf{h}^T \mathbf{z} \\ &= \sup_{\lambda, \eta \geq 0, \nu} \eta J_\infty \left(\frac{\mathbf{h}}{\eta \sqrt{n}} - \left(\frac{\nu}{\eta} - \theta \right) \tilde{\boldsymbol{\mu}}; \frac{\lambda}{\eta} \right) - \frac{\nu^2}{2\eta} - \frac{1}{2\eta n} \|\mathbf{h}\|_{\ell_2}^2 + \frac{\nu}{\eta \sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} - \frac{\eta}{2} \alpha^2 - \frac{\lambda}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_1}}, \end{aligned}$$

where $J_\infty(\mathbf{x}; \lambda)$ is defined as

$$(D.17) \quad J_\infty(\mathbf{x}; \lambda) = \min_{\mathbf{v}} \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_{\ell_2}^2 + \lambda \|\mathbf{v}\|_{\ell_\infty}.$$

Our next step is scalarization of the optimization (D.16) in the large sample limit (as $n \rightarrow \infty$), and a challenge along this way is that the function $J(\mathbf{x}; \lambda)$ is not a separable function over the entries of \mathbf{x} . To cope with this problem, we propose an alternative representation of this function that involves an additional variable t_0 .

We write

$$(D.18) \quad \begin{aligned} J_\infty(\mathbf{x}; \lambda) &= \min_{\mathbf{v}} \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_{\ell_2}^2 + \lambda \|\mathbf{v}\|_{\ell_\infty} \\ &= \min_{\mathbf{v}, t \geq 0} \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_{\ell_2}^2 + \lambda t \quad \text{subject to} \quad \|\mathbf{v}\|_{\ell_\infty} \leq t \\ &= \min_{t \geq 0} \frac{1}{2} \|\text{ST}(\mathbf{x}; t)\|_{\ell_2}^2 + \lambda t. \end{aligned}$$

Let $t_0 = \|\boldsymbol{\mu}\|_{\ell_1} t$ and $\lambda_0 = \frac{\lambda}{\|\boldsymbol{\mu}\|_{\ell_1}}$. Similar to the trick of ‘artificial’ boundedness that we used in applying the CGMT framework (e.g., cf. explanation after (B.16) and Appendix E.3), we continue by the ansatz that the optimal value of λ_0 and t_0 remain bounded as $n \rightarrow \infty$. After we take the limit of the Lagrangian to obtain a scalar auxiliary optimization (AO) problem, this ansatz is verified by the boundedness of solutions of the AO problem.

For $\mathbf{x} = \frac{c_0}{\sqrt{n}} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}}$ we have

$$(D.19) \quad \begin{aligned} &\lim_{n \rightarrow \infty} \frac{1}{2} \|\text{ST}(\mathbf{x}; t)\|_{\ell_2}^2 + \lambda t \\ &= \lim_{n \rightarrow \infty} \frac{1}{2} \sum_{i=1}^d \text{ST}\left(x_i; \frac{t_0}{\|\boldsymbol{\mu}\|_{\ell_1}}\right)^2 + \lambda_0 t_0 \\ &= \lim_{n \rightarrow \infty} \frac{1}{2} \sum_{i=1}^d \text{ST}\left(\frac{c_0}{\sqrt{n}} h_i - c_1 \tilde{\mu}_i; \frac{t_0}{\|\boldsymbol{\mu}\|_{\ell_1}}\right)^2 + \lambda_0 t_0 \\ &\stackrel{(a)}{=} \lim_{n \rightarrow \infty} \frac{1}{2d} \sum_{i=1}^d \text{ST}\left(\frac{c_0}{\sqrt{\delta}} h_i - c_1 \frac{\sqrt{d} \mu_i}{\sigma_{M,2}}; \frac{t_0}{\sigma_{M,1}}\right)^2 + \lambda_0 t_0 \\ &= \frac{1}{2} \mathbb{E} \left[\text{ST}\left(\frac{c_0}{\sqrt{\delta}} h - c_1 \frac{M}{\sigma_{M,2}}; \frac{t_0}{\sigma_{M,1}}\right)^2 \right] + \lambda_0 t_0, \\ &= f(c_0, c_1; t_0) + \lambda_0 t_0, \end{aligned}$$

with high probability. In (a) we used the fact that as $n \rightarrow \infty$, we have $\|\boldsymbol{\mu}\|_{\ell_1} \rightarrow \sqrt{d} \sigma_{M,1}$ along with the identity $\frac{1}{a^2} \text{ST}(a\mathbf{x}; a\lambda) = \text{ST}(\mathbf{x}; \lambda)$.

Note that this is a pointwise convergence. However, the left hand side is a convex function of t and it's minimizer satisfies $t \leq \|\mathbf{x}\|_\infty \leq c_0 + c_1$ and so belongs to a compact set. Therefore, by applying the convexity lemma, see e.g, [39, Lemma 7.75], [63, Lemma B1], we can change the order of limit and minimization, and get that for $\mathbf{x} = \frac{c_0}{\sqrt{n}}\mathbf{h} - c_1\tilde{\boldsymbol{\mu}}$,

$$\begin{aligned}
 & \lim_{n \rightarrow \infty} J_\infty(\mathbf{x}; \lambda) \\
 &= \lim_{n \rightarrow \infty} \min_{t \geq 0} \left\{ \frac{1}{2} \|\text{ST}(\mathbf{x}; t)\|_{\ell_2}^2 + \lambda t \right\} \\
 (D.20) \quad &= \min_{t_0 \geq 0} \{f(c_0, c_1; t_0) + \lambda_0 t_0\},
 \end{aligned}$$

in probability. In addition, as $n \rightarrow \infty$ we have

$$(D.21) \quad \frac{1}{n} \|\mathbf{h}\|_{\ell_2}^2 \rightarrow \frac{1}{\delta}, \quad \frac{1}{\sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} \rightarrow 0, \quad \|\boldsymbol{\mu}\|_{\ell_1} \rightarrow \sigma_{M,1} \sqrt{d},$$

with high probability.

Using the above limits, we see that the objective function (D.16) converges pointwise to the following function:

$$(D.22) \quad \min_{t_0 \geq 0} \eta \left\{ f\left(\frac{1}{\eta}, \frac{\nu}{\eta} - \theta; t_0\right) + \frac{\lambda_0}{\eta} t_0 \right\} - \frac{\nu^2}{2\eta} - \frac{1}{2\eta\delta} - \frac{\eta}{2} \alpha^2 - \frac{\lambda_0}{\varepsilon_0}$$

Note that (E.4) is the dual optimization and hence is a concave problem. We apply the convexity lemma [63, Lemma B.2] to conclude that the objective value in (E.4) also converges to the supremum of function (E.8) over $\lambda_0, \eta \geq 0, \nu$. Therefore the solution of optimization (D.16) converges to the solution of the following optimization problem:

$$(D.23) \quad \sup_{\lambda_0, \eta \geq 0, \nu} \min_{t_0 \geq 0} \eta f\left(\frac{1}{\eta}, \frac{\nu}{\eta} - \theta; t_0\right) + \lambda_0 t_0 - \frac{\nu^2}{2\eta} - \frac{1}{2\eta\delta} - \frac{\eta}{2} \alpha^2 - \frac{\lambda_0}{\varepsilon_0}$$

Note that the above objective is linear in λ_0 . Therefore the optimal t_0^* should satisfy $t_0^* \leq \frac{1}{\varepsilon_0}$. Otherwise $\lambda_0^* = \infty$ which makes the above max-min value unbounded, and this is a contradiction because the above problem involves minimization over t_0 and it is easy to see that by choosing $t_0 = 0$ the optimal objective value over $\{\lambda_0, \eta \geq 0, \nu\}$ becomes zero.

Therefore, we can assume $t_0 \leq \frac{1}{\varepsilon_0}$ which yields $\lambda_0^*(t_0 - \frac{1}{\varepsilon_0}) = 0$. This simplifies the problem (D.23) to

$$(D.24) \quad \sup_{\eta \geq 0, \nu} \min_{0 \leq t_0 \leq \frac{1}{\varepsilon_0}} \eta f\left(\frac{1}{\eta}, \frac{\nu}{\eta} - \theta; t_0\right) - \frac{\nu^2}{2\eta} - \frac{1}{2\eta\delta} - \frac{\eta}{2} \alpha^2$$

Since $f(c_0, c_1; t)$ is decreasing in t , the optimal value t_0^* is given by $t_0^* = \frac{1}{\varepsilon_0}$ and the problem is further simplified and along with (D.16) implies that

$$(D.25) \quad \begin{aligned} & \min_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})} -\frac{1}{\sqrt{n}} \mathbf{h}^T \mathbf{z} \\ &= \sup_{\eta \geq 0, \nu} \eta f\left(\frac{1}{\eta}, \frac{\nu}{\eta} - \theta; \frac{1}{\varepsilon_0}\right) - \frac{\nu^2}{2\eta} - \frac{1}{2\eta\delta} - \frac{\eta}{2} \alpha^2 \end{aligned}$$

Now similar to the proof of Lemma 4.2 we use Equation (E.3) to write

$$(D.26) \quad \begin{aligned} & \omega(\alpha, \theta, \varepsilon_0) \\ &= \lim_{n \rightarrow \infty} \omega_s(\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})) \\ &= - \sup_{\eta \geq 0, \nu} \sqrt{\delta} \left\{ \eta f\left(\frac{1}{\eta}, \frac{\nu}{\eta} - \theta; \frac{1}{\varepsilon_0}\right) - \frac{\nu^2}{2\eta} - \frac{1}{2\eta\delta} - \frac{\eta}{2} \alpha^2 \right\} \\ &= \min_{\eta \geq 0, \nu} \sqrt{\delta} \left\{ \frac{\nu^2}{2\eta} + \frac{1}{2\eta\delta} + \frac{\eta}{2} \alpha^2 - \eta f\left(\frac{1}{\eta}, \frac{\nu}{\eta} - \theta; \frac{1}{\varepsilon_0}\right) \right\}. \end{aligned}$$

This completes the proof. ■

Proof [Part (b)] The proof of this parts proceeds along the same lines of Theorem 4.3 for the special case of $p = 1$, $q = \infty$. However, it requires a slightly different treatment as in part (a) because the function $\|\cdot\|_{\ell_q}^q$ and therefore J_q given by (4.5) are not well defined in this case.

Here we only highlight the modifications that are needed to the proof of Theorem 4.3 to apply it for case of $q = \infty$.

We proceed the exact same derivation that yields (B.22), repeated here for convenience:

$$(D.27) \quad \begin{aligned} & \min_{\boldsymbol{\theta}, \boldsymbol{\mu}, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \lambda, \eta \geq 0, \tilde{\eta}} \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \left(\|\boldsymbol{\theta}\|_{\ell_q} - \gamma \right) + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_{\boldsymbol{\mu}}^\perp \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta} \\ &+ \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon\gamma - \boldsymbol{\theta}^T \boldsymbol{\mu})_{\ell_2} + \alpha g \right)_+^2 \right]} + \eta \left(\|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\|_{\ell_2} - \alpha \right) + \tilde{\eta} (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} - \theta) \end{aligned}$$

We substitute for $\|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\|_{\ell_2}$ using the identity $\|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\|_{\ell_2} = \min_{\tau \geq 0} \frac{\|\boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\theta}\|_{\ell_2}^2}{2\tau} + \frac{\tau}{2}$

to get

$$\begin{aligned}
& \min_{\boldsymbol{\theta}, \boldsymbol{\mu}, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \lambda, \eta \geq 0, \tilde{\eta}} \min_{\tau \geq 0} \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \|\boldsymbol{\theta}\|_{\ell_\infty} - 2\lambda\gamma + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_\mu^\perp \Sigma^{1/2} \boldsymbol{\theta} \\
& + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon\gamma - \theta \|\boldsymbol{\mu}\|_{\ell_2}) + \alpha g \right)_+^2 \right]} \\
& + \frac{\eta}{2\tau} \left\| \Sigma^{\frac{1}{2}} \boldsymbol{\theta} \right\|_{\ell_2}^2 + \frac{\eta\tau}{2} - \eta\alpha + \tilde{\eta} (\tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} - \theta). \tag{D.28}
\end{aligned}$$

Specializing it to $\Sigma = \mathbf{I}$ and $q = \infty$, the optimization over $\boldsymbol{\theta}$ takes the form

$$\min_{\boldsymbol{\theta}} \left(1 + \frac{\eta}{2\tau} \right) \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \|\boldsymbol{\theta}\|_{\ell_\infty} + \frac{\beta}{\sqrt{n}} \mathbf{h}^T \mathbf{P}_\mu^\perp \boldsymbol{\theta} + \tilde{\eta} \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta}. \tag{D.29}$$

Note that

$$\begin{aligned}
\frac{\beta}{\sqrt{n}} \mathbf{P}_\mu^\perp \mathbf{h} - \frac{\tilde{\eta} \boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_{\ell_2}} &= \frac{\beta}{\sqrt{n}} \mathbf{h} - \frac{\tilde{\eta} \boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_{\ell_2}} - \frac{\beta}{\sqrt{n}} \mathbf{P}_\mu \mathbf{h} \\
&= \frac{\beta}{\sqrt{n}} \mathbf{h} - \left(\frac{\tilde{\eta}}{\|\boldsymbol{\mu}\|_{\ell_2}} + \frac{\beta}{\sqrt{n} \|\boldsymbol{\mu}\|_{\ell_2}} \mathbf{h}^T \tilde{\boldsymbol{\mu}} \right) \boldsymbol{\mu},
\end{aligned}$$

where $\tilde{\boldsymbol{\mu}}^T \mathbf{h} \sim \mathcal{N}(0, 1)$ since $\|\tilde{\boldsymbol{\mu}}\|_{\ell_2} = 1$, and $\|\boldsymbol{\mu}\|_{\ell_2} \rightarrow \sigma_{M,2}$ which implies that the last term in the right-hand side is dominated by the second term. Therefore in the asymptotic regime $n \rightarrow \infty$, we can equivalently work replace $\mathbf{P}_\mu^\perp \mathbf{h}$ by \mathbf{h} and by using the symmetry of the Gaussian distribution work with

$$\min_{\boldsymbol{\theta}} \left(1 + \frac{\eta}{2\tau} \right) \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \|\boldsymbol{\theta}\|_{\ell_\infty} - \frac{\beta}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\theta} + \tilde{\eta} \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta}. \tag{D.30}$$

We let $\mathbf{x} = \frac{\beta}{\sqrt{n}} \mathbf{h} - \tilde{\eta} \tilde{\boldsymbol{\mu}}$ and consider the change of variable $\mathbf{u} := 2(1 + \frac{\eta}{2\tau}) \boldsymbol{\theta}$ and write

$$\begin{aligned}
& \min_{\boldsymbol{\theta}} \left(1 + \frac{\eta}{2\tau} \right) \|\boldsymbol{\theta}\|_{\ell_2}^2 + 2\lambda \|\boldsymbol{\theta}\|_{\ell_\infty} - \frac{\beta}{\sqrt{n}} \mathbf{h}^T \boldsymbol{\theta} + \tilde{\eta} \tilde{\boldsymbol{\mu}}^T \boldsymbol{\theta} \\
&= \min_{\mathbf{u}} \frac{1}{2} \left(1 + \frac{\eta}{2\tau} \right)^{-1} \left\{ \frac{1}{2} \|\mathbf{u}\|_{\ell_2}^2 + 2\lambda \|\mathbf{u}\|_{\ell_\infty} - \mathbf{x}^T \mathbf{u} \right\} \\
&= \min_{\mathbf{u}} \frac{1}{2} \left(1 + \frac{\eta}{2\tau} \right)^{-1} \left\{ \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_{\ell_2}^2 + 2\lambda \|\mathbf{u}\|_{\ell_\infty} - \frac{\|\mathbf{x}\|_{\ell_2}^2}{2} \right\} \\
&= \frac{1}{2} \left(1 + \frac{\eta}{2\tau} \right)^{-1} J_\infty(\mathbf{x}; 2\lambda) - \frac{1}{4} \left(1 + \frac{\eta}{2\tau} \right)^{-1} \|\mathbf{x}\|_{\ell_2}^2 \tag{D.31}
\end{aligned}$$

Using (D.31) and substituting for \mathbf{x} in (D.27), our AO problem becomes

$$(D.32) \quad \min_{\theta, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \lambda, \eta \geq 0, \tilde{\eta}} \min_{\tau \geq 0} \frac{1}{2} \left(1 + \frac{\eta}{2\tau}\right)^{-1} J_{\infty} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} - \tilde{\eta} \tilde{\boldsymbol{\mu}}; 2\lambda \right) - \frac{1}{4} \left(1 + \frac{\eta}{2\tau}\right)^{-1} \left\| \frac{\beta}{\sqrt{n}} \mathbf{h} - \tilde{\eta} \tilde{\boldsymbol{\mu}} \right\|_{\ell_2}^2$$

$$- 2\lambda\gamma + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon\gamma - \theta \|\boldsymbol{\mu}\|_{\ell_2}) + \alpha g \right)_+^2 \right]} + \frac{\eta\tau}{2} - \eta\alpha - \tilde{\eta}\theta$$

Our next step is to scalarize the AO problem by taking the asymptotic limit of the objective.

We have

$$\lim_{n \rightarrow \infty} \left\| \frac{\beta}{\sqrt{n}} \mathbf{h} - \tilde{\eta} \tilde{\boldsymbol{\mu}} \right\|_{\ell_2}^2 = \frac{\beta^2}{\delta} + \tilde{\eta}^2, \quad \lim_{n \rightarrow \infty} \|\boldsymbol{\mu}\|_{\ell_2} = \sigma_{M,2}.$$

in probability. Also by using (D.20) we have

$$\lim_{n \rightarrow \infty} J_{\infty} \left(\frac{\beta}{\sqrt{n}} \mathbf{h} - \tilde{\eta} \tilde{\boldsymbol{\mu}}; 2\lambda \right) = \min_{t_0 \geq 0} \{ f(\beta, \tilde{\eta}; t_0) + 2\lambda_0 t_0 \},$$

with $\lambda_0 = \frac{\lambda}{\|\boldsymbol{\mu}\|_{\ell_1}}$. Using these limits in the AO problem (D.32) we obtain the following scalar AO problem

$$(D.33) \quad \min_{\theta, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \lambda, \eta \geq 0, \tilde{\eta}} \min_{\tau, t_0 \geq 0} \frac{1}{2} \left(1 + \frac{\eta}{2\tau}\right)^{-1} (f(\beta, \tilde{\eta}; t_0) + 2\lambda_0 t_0) - \frac{1}{4} \left(1 + \frac{\eta}{2\tau}\right)^{-1} \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right)$$

$$- 2\lambda_0 \gamma_0 + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right]} + \frac{\eta\tau}{2} - \eta\alpha - \tilde{\eta}\theta,$$

where we recall our notation $\varepsilon_0 = \frac{\varepsilon}{\|\boldsymbol{\mu}\|_{\ell_1}}$ and $\gamma_0 = \gamma \|\boldsymbol{\mu}\|_{\ell_1}$.

Now note that objective function (D.33) is linear in λ_0 and therefore the optimal t_0^* should satisfy $t_0^* \leq 2\gamma_0 \left(1 + \frac{\eta}{2\tau}\right)$, otherwise $\lambda_0^* = \infty$ which makes the above max-min value unbounded. As such, we also have $\lambda_0^* = 0$ which further simplifies the problem as follows:

$$(D.34) \quad \min_{\theta, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \eta \geq 0, \tilde{\eta}} \min_{0 \leq \tau, 0 \leq t_0 \leq 2\gamma_0(1 + \frac{\eta}{2\tau})} \frac{1}{2} \left(1 + \frac{\eta}{2\tau}\right)^{-1} f(\beta, \tilde{\eta}; t_0) - \frac{1}{4} \left(1 + \frac{\eta}{2\tau}\right)^{-1} \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2 \right)$$

$$+ \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right]} + \frac{\eta\tau}{2} - \eta\alpha - \tilde{\eta}\theta,$$

Since $f(c_0, c_1; t_0)$ is decreasing in t_0 , the optimal value of t_0 is given by $t_0^* = 2\gamma_0(1 + \frac{\eta}{2\tau})$ which results in the following AO problem:

$$(D.35) \quad \min_{\theta, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \eta \geq 0, \tilde{\eta}} \min_{\tau \geq 0} \frac{1}{2} \left(1 + \frac{\eta}{2\tau}\right)^{-1} f\left(\beta, \tilde{\eta}; 2\gamma_0\left(1 + \frac{\eta}{2\tau}\right)\right) - \frac{1}{4} \left(1 + \frac{\eta}{2\tau}\right)^{-1} \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2\right) \\ + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right]} + \frac{\eta\tau}{2} - \eta\alpha - \tilde{\eta}\theta.$$

Our final step of simplification is to solve for τ . To this end, we define the function

$$R\left(\frac{\eta}{\tau}\right) := \frac{1}{2} \left(1 + \frac{\eta}{2\tau}\right)^{-1} f\left(\beta, \tilde{\eta}; 2\gamma_0\left(1 + \frac{\eta}{2\tau}\right)\right) - \frac{1}{4} \left(1 + \frac{\eta}{2\tau}\right)^{-1} \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2\right),$$

where we make the dependence on $\frac{\eta}{\tau}$ explicit in the notation. Setting derivative of the AO objective with respect to η , to zero we obtain

$$\frac{1}{\tau} R'\left(\frac{\eta}{\tau}\right) + \frac{\tau}{2} - \alpha = 0.$$

Setting derivative with respect to τ to zero gives

$$-\frac{\eta}{\tau^2} R'\left(\frac{\eta}{\tau}\right) + \frac{\eta}{2} = 0.$$

Combining the above two optimality condition implies that $\eta(1 - \frac{\alpha}{\tau}) = 0$. So either $\alpha = \tau$ or $\eta = 0$. If $\eta = 0$, then it is clear that the terms involving τ in the AO problem would vanish and therefore the value of τ does not matter. So in this case, we can as well assume $\tau = \alpha$. Substituting for τ the AO problem further simplifies to

$$(D.36) \quad \min_{\theta, \gamma \geq 0, \alpha \geq 0} \max_{\beta, \eta \geq 0, \tilde{\eta}} \frac{1}{2} \left(1 + \frac{\eta}{2\alpha}\right)^{-1} f\left(\beta, \tilde{\eta}; 2\gamma_0\left(1 + \frac{\eta}{2\alpha}\right)\right) - \frac{1}{4} \left(1 + \frac{\eta}{2\alpha}\right)^{-1} \left(\frac{\beta^2}{\delta} + \tilde{\eta}^2\right) \\ + \beta \sqrt{\mathbb{E} \left[\left((1 + \varepsilon_0 \gamma_0 - \theta \sigma_{M,2}) + \alpha g \right)_+^2 \right]} - \frac{\eta\alpha}{2} - \tilde{\eta}\theta.$$

This concludes the proof of part (b). ■

Proof [Part (c)] The proof of part (c) follows along the same lines of the proof of Theorem 6.4 (and Theorem 4.5 for isotropic case). But similar to previous parts, we need to make slight modifications to the proof.

Note that in our derivation of the AO problem (6.9), we replaced the constraint $\|\mathbf{u}\|_{\ell_q} \leq \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}}$ with the equivalent constraint $\|\mathbf{u}\|_{\ell_q}^q \leq (\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}})^q$, see (E.11) for more details. The benefit of this alternative representation is that it results in the Moreau-envelope $e_{q,\Sigma}$ of the $\|\cdot\|_{\ell_q}^q$ function, see (E.15), which is separable over the samples. As a result, in the isotropic case the expected Moreau envelope reduces to the expected of the one-dimensional function J_q , given by (4.5), (C.9).

However, for $q = \infty$ the function $\|\cdot\|_{\ell_q}^q$ is not well-defined and requires a slightly different treatment. In this case we stay with the original constraint $\|\mathbf{u}\|_{\ell_q} \leq \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}}$. Proceeding along the same derivations of AO problem (4.16), it is straightforward to see that this results in the following AO problem for the non-separable regime:

$$\begin{aligned} & \max_{0 \leq \beta, \tau_h} \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) \\ D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) &= \frac{\beta \tau_g}{2} + L\left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2}\theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta}\right) \\ & - \min_{\lambda_0 \geq 0, \nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \frac{\gamma_0 \tau_h}{\alpha} + \frac{\nu^2}{2} - \tilde{\mathbb{E}}\left(\beta, \frac{\tau_h \theta}{\alpha} + \nu; \lambda_0\right) \right\} + \frac{\alpha \tau_h}{2} \right] \end{aligned} \quad (\text{D.37})$$

with

$$(\text{D.38}) \quad \tilde{\mathbb{E}}(c_0, c_1; \lambda_0) := \lim_{n \rightarrow \infty} J_{\infty} \left(\frac{c_0}{\sqrt{n}} \mathbf{h} - c_1 \tilde{\boldsymbol{\mu}}; \lambda_0 \|\boldsymbol{\mu}\|_{\ell_1} \right),$$

and $J_{\infty}(\mathbf{x}; \lambda)$ given by (D.17). Using (D.20), we have

$$(\text{D.39}) \quad \tilde{\mathbb{E}}(c_0, c_1; \lambda_0) = \min_{t_0 \geq 0} \{f(c_0, c_1; t_0) + \lambda_0 t_0\},$$

Substituting for $\tilde{\mathbb{E}}$ function in the AO problem (D.37) results in

$$\begin{aligned} & \max_{0 \leq \beta, \tau_h} \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) \\ D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) &= \frac{\beta \tau_g}{2} + L\left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2}\theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta}\right) \\ & - \min_{\lambda_0 \geq 0, \nu} \sup_{t_0 \geq 0} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \lambda_0 \left(\frac{\gamma_0 \tau_h}{\alpha} - t_0 \right) + \frac{\nu^2}{2} - f\left(\beta, \frac{\tau_h \theta}{\alpha} + \nu; t_0\right) \right\} + \frac{\alpha \tau_h}{2} \right] \end{aligned} \quad (\text{D.40})$$

Note that the above objective is linear in λ_0 . Clearly, the optimal value t_0^* should satisfy $t_0^* \leq \frac{\gamma_0 \tau_h}{\alpha}$; otherwise $\lambda_0^* = \infty$ which makes the objective value

unbounded. For $t_0 \leq \frac{\gamma_0 \tau_h}{\alpha}$, we have $\lambda_0^* = 0$. Therefore, t_0 only appears in the term $f(\beta, \frac{\tau_h \theta}{\alpha} + \nu; t_0)$. Given that $f(c_0, c_1; t_0)$ is decreasing in t_0 , we have $t_0^* = \frac{\gamma_0 \tau_h}{\alpha}$.

We substitute for t_0^* in the AO problem to obtain

$$\begin{aligned}
 & \max_{0 \leq \beta, \tau_h} \min_{\theta, 0 \leq \alpha, \gamma_0, \tau_g} D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) \\
 & D_{\text{ns}}(\alpha, \gamma_0, \theta, \tau_g, \beta, \tau_h) = \frac{\beta \tau_g}{2} + L \left(\sqrt{\alpha^2 + \theta^2}, \sigma_{M,2} \theta - \varepsilon_0 \gamma_0, \frac{\tau_g}{\beta} \right) \\
 & \quad - \min_{\nu} \left[\frac{\alpha}{\tau_h} \left\{ \frac{\beta^2}{2\delta} + \frac{\nu^2}{2} - f \left(\beta, \frac{\tau_h \theta}{\alpha} + \nu; \frac{\gamma_0 \tau_h}{\alpha} \right) \right\} + \frac{\alpha \tau_h}{2} \right].
 \end{aligned}
 \tag{D.41}$$

This completes the proof of part (c). \blacksquare

APPENDIX E: PROOF OF TECHNICAL LEMMAS

E.1. Proof of Lemma 2.1 By definition, we have

$$\begin{aligned}
 \text{SA}(\widehat{\boldsymbol{\theta}}) &:= \mathbb{E}[\mathbb{1}(\hat{y} = y)] = \mathbb{P}(y \langle \mathbf{x}, \widehat{\boldsymbol{\theta}} \rangle > 0) \\
 &= \mathbb{P} \left(y \langle y \boldsymbol{\mu} + \Sigma^{1/2} \mathbf{z}, \widehat{\boldsymbol{\theta}} \rangle > 0 \right) \\
 &= \mathbb{P} \left(\langle \boldsymbol{\mu} + \Sigma^{1/2} \mathbf{z}, \widehat{\boldsymbol{\theta}} \rangle > 0 \right) \\
 &= \mathbb{P} \left(\langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}} \rangle + \left\| \Sigma^{1/2} \widehat{\boldsymbol{\theta}} \right\|_{\ell_2} Z > 0 \right) \\
 &= \Phi \left(\frac{\langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}} \rangle}{\left\| \Sigma^{1/2} \widehat{\boldsymbol{\theta}} \right\|_{\ell_2}} \right),
 \end{aligned}
 \tag{E.1}$$

where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_d)$ and $Z \sim \mathcal{N}(0, 1)$.

Likewise for the adversarial risk we have

$$\begin{aligned}
\text{RA}(\widehat{\boldsymbol{\theta}}) &:= \mathbb{E} \left[\min_{\|\boldsymbol{\delta}\|_{\ell_p} \leq \varepsilon} \mathbb{1}(y\langle \mathbf{x} + \boldsymbol{\delta}, \widehat{\boldsymbol{\theta}} \rangle \geq 0) \right] \\
&\stackrel{(a)}{=} \mathbb{E} \left[\mathbb{1}(y\langle \mathbf{x}, \widehat{\boldsymbol{\theta}} \rangle - \varepsilon \|\widehat{\boldsymbol{\theta}}\|_{\ell_q} \geq 0) \right] \\
&= \mathbb{P} \left(y\langle \mathbf{x}, \widehat{\boldsymbol{\theta}} \rangle - \varepsilon \|\widehat{\boldsymbol{\theta}}\|_{\ell_q} \geq 0 \right) \\
&= \mathbb{P} \left(y\langle y\boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\mathbf{z}, \widehat{\boldsymbol{\theta}} \rangle - \varepsilon \|\widehat{\boldsymbol{\theta}}\|_{\ell_q} \geq 0 \right) \\
&\stackrel{(b)}{=} \mathbb{P} \left(\langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}} \rangle + \|\boldsymbol{\Sigma}^{1/2}\widehat{\boldsymbol{\theta}}\|_{\ell_2} Z - \varepsilon \|\widehat{\boldsymbol{\theta}}\|_{\ell_q} \geq 0 \right) \\
&= \Phi \left(\frac{\langle \boldsymbol{\mu}, \widehat{\boldsymbol{\theta}} \rangle - \varepsilon \|\widehat{\boldsymbol{\theta}}\|_{\ell_q}}{\|\boldsymbol{\Sigma}^{1/2}\widehat{\boldsymbol{\theta}}\|_{\ell_2}} \right),
\end{aligned} \tag{E.2}$$

where (a) we used that $\langle \boldsymbol{\delta}, \widehat{\boldsymbol{\theta}} \rangle \geq -\|\boldsymbol{\delta}\|_{\ell_p} \|\widehat{\boldsymbol{\theta}}\|_{\ell_q} \geq -\varepsilon \|\widehat{\boldsymbol{\theta}}\|_{\ell_q}$, using Hölder inequality (with $\frac{1}{p} + \frac{1}{q} = 1$) and that $\|\boldsymbol{\delta}\|_{\ell_p} \leq \varepsilon$, with equality achieving for some $\boldsymbol{\delta}$ in this set. In (b), we used the symmetry of Gaussian distribution.

E.2. Proof of Lemma 4.2 We first note that

$$\min_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})} -\frac{1}{\sqrt{n}} \mathbf{h}^T \mathbf{z} = -\frac{1}{\sqrt{n/d}} \sup_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} \rightarrow -\frac{1}{\sqrt{\delta}} \omega(\alpha, \theta, \varepsilon_0), \tag{E.3}$$

in probability, using the fact that \mathbf{h}/\sqrt{d} is asymptotically uniform on the unit sphere, and for $\mathcal{S} \in \mathbb{S}^{d-1}$ the function $f(\mathbf{u}) = \sup_{\mathbf{z} \in \mathcal{S}} \mathbf{z}^T \mathbf{u}$ is Lipschitz. Therefore, using the concentration of Lipschitz functions on the sphere (see e.g. [66, Theorem 5.2.2]), $f(\mathbf{u})$ concentrates around its mean $\mathbb{E} f(\mathbf{u}) = \omega_s(\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu}))$. More precisely,

$$\mathbb{P} \left\{ \left| \sup_{\mathbf{z} \in \mathcal{S}} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} - \omega_s(\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})) \right| \right\} \leq 2e^{-cdt^2},$$

for an absolute constant $c > 0$ and for every $t \geq 0$. Therefore, by invoking the assumption on the convergence of spherical width, cf. Assumption 3, we arrive at

$$\lim_{d \rightarrow \infty} \mathbb{P} \left\{ \left| \sup_{\mathbf{z} \in \mathcal{S}} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} - \omega(\alpha, \theta, \varepsilon_0) \right| \geq \eta \right\} = 0, \quad \forall \eta > 0.$$

Therefore, $\sup_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0)} \frac{1}{\sqrt{d}} \mathbf{h}^T \mathbf{z} \rightarrow \omega(\alpha, \theta, \varepsilon_0)$, in probability.

To evaluate the left hand side, we form the Lagrangian corresponding to the set \mathcal{S} . Let $\tilde{\boldsymbol{\mu}} := \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_{\ell_2}}$ and consider the change of variable $\mathbf{u} := \mathbf{z} + \theta \tilde{\boldsymbol{\mu}}$. We then have

$$\begin{aligned}
& \min_{\mathbf{z} \in \mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})} -\frac{1}{\sqrt{n}} \mathbf{h}^T \mathbf{z} \\
&= \sup_{\lambda, \eta \geq 0, \nu} \min_{\mathbf{u}} -\frac{1}{\sqrt{n}} \mathbf{h}^T (\mathbf{u} - \theta \tilde{\boldsymbol{\mu}}) + \lambda \left(\|\mathbf{u}\|_{\ell_q}^q - \left(\frac{1}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}} \right)^q \right) + \frac{\eta}{2} (\|\mathbf{u} - \theta \tilde{\boldsymbol{\mu}}\|_{\ell_2}^2 - \alpha^2) + \nu \tilde{\boldsymbol{\mu}}^T (\mathbf{u} - \theta \tilde{\boldsymbol{\mu}}) \\
&= \sup_{\lambda, \eta \geq 0, \nu} \min_{\mathbf{u}} \frac{\eta}{2} \left\| \mathbf{u} - \theta \tilde{\boldsymbol{\mu}} + \frac{\nu}{\eta} \tilde{\boldsymbol{\mu}} - \frac{\mathbf{h}}{\eta \sqrt{n}} \right\|_{\ell_2}^2 + \lambda \left\| \mathbf{u} \right\|_{\ell_q}^q - \frac{1}{2\eta} \left\| \nu \tilde{\boldsymbol{\mu}} - \frac{\mathbf{h}}{\sqrt{n}} \right\|_{\ell_2}^2 - \frac{\eta}{2} \alpha^2 - \frac{\lambda}{(\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p})^q} \\
&= \sup_{\lambda, \eta \geq 0, \nu} \min_{\mathbf{u}} \eta \left[\frac{1}{2} \left\| \mathbf{u} + \left(\frac{\nu}{\eta} - \theta \right) \tilde{\boldsymbol{\mu}} - \frac{\mathbf{h}}{\eta \sqrt{n}} \right\|_{\ell_2}^2 + \frac{\lambda}{\eta} \|\mathbf{u}\|_{\ell_q}^q \right] - \frac{1}{2\eta} \left\| \nu \tilde{\boldsymbol{\mu}} - \frac{\mathbf{h}}{\sqrt{n}} \right\|_{\ell_2}^2 - \frac{\eta}{2} \alpha^2 - \frac{\lambda}{(\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p})^q} \\
&= \sup_{\lambda, \eta \geq 0, \nu} \eta \sum_{i=1}^d J_q \left(\frac{h_i}{\eta \sqrt{n}} - \left(\frac{\nu}{\eta} - \theta \right) \tilde{\mu}_i; \frac{\lambda}{\eta} \right) - \frac{1}{2\eta} \left\| \nu \tilde{\boldsymbol{\mu}} - \frac{\mathbf{h}}{\sqrt{n}} \right\|_{\ell_2}^2 - \frac{\eta}{2} \alpha^2 - \frac{\lambda}{(\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p})^q} \\
&\quad \text{(E.4)} \\
&= \sup_{\lambda, \eta \geq 0, \nu} \eta \sum_{i=1}^d J_q \left(\frac{h_i}{\eta \sqrt{n}} - \left(\frac{\nu}{\eta} - \theta \right) \tilde{\mu}_i; \frac{\lambda}{\eta} \right) - \frac{\nu^2}{2\eta} - \frac{1}{2\eta n} \|\mathbf{h}\|_{\ell_2}^2 + \frac{\nu}{\eta \sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} - \frac{\eta}{2} \alpha^2 - \frac{\lambda}{(\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p})^q}
\end{aligned}$$

Recall that $\mathbf{h} \sim \mathcal{N}(0, \mathbf{I}_d)$. As $n \rightarrow \infty$ and $n/d \rightarrow \delta$, we have

$$\text{(E.5)} \quad \frac{1}{n} \|\mathbf{h}\|_{\ell_2}^2 \rightarrow \frac{1}{\delta}, \quad \frac{1}{\sqrt{n}} \tilde{\boldsymbol{\mu}}^T \mathbf{h} \rightarrow 0,$$

in probability. In addition,

$$\text{(E.6)} \quad \|\boldsymbol{\mu}\|_{\ell_p} \rightarrow \sigma_{M,p} d^{\frac{1}{p} - \frac{1}{2}} = \sigma_{M,p} d^{\frac{1}{2} - \frac{1}{q}},$$

in probability. Using the identity $J_q(x; \lambda) = c^2 J_q(x/c; \lambda c^{q-2})$ and letting $\lambda_0 := \lambda d^{1-\frac{q}{2}}$ we have

$$\begin{aligned}
\sum_{i=1}^d J_q \left(\frac{h_i}{\eta \sqrt{n}} - \left(\frac{\nu}{\eta} - \theta \right) \tilde{\mu}_i; \frac{\lambda}{\eta} \right) &= \frac{1}{d} \sum_{i=1}^d J_q \left(\frac{\sqrt{d} h_i}{\eta \sqrt{n}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{\sqrt{d} \mu_i}{\sigma_{M,2}}; \frac{\lambda}{\eta} d^{1-q/2} \right) \\
&= \frac{1}{d} \sum_{i=1}^d J_q \left(\frac{h_i}{\eta \sqrt{\delta}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{\sqrt{d} \mu_i}{\sigma_{M,2}}; \frac{\lambda_0}{\eta} \right).
\end{aligned}$$

Since $J_q(x; \lambda) \leq \frac{1}{2} x^2$, the function J_q is pseudo-lipschitz of order 2 and by an

application of [3, Lemma 5], we have

(E.7)

$$\lim_{n \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d J_q \left(\frac{h_i}{\eta \sqrt{\delta}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{\sqrt{d} \mu_i}{\sigma_{M,2}}; \frac{\lambda_0}{\eta} \right) = \mathbb{E} \left[J_q \left(\frac{h}{\eta \sqrt{\delta}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{M}{\sigma_{M,2}}; \frac{\lambda_0}{\eta} \right) \right],$$

almost surely, where the expectation in the last line is taken with respect to the independent random variables $h \sim \mathcal{N}(0, 1)$ and $M \sim \mathbb{P}_M$.

Using the above limits, we see that the objective function (E.4) converges pointwise to the following function:

$$(E.8) \quad \eta \mathbb{E} \left[J_q \left(\frac{h}{\eta \sqrt{\delta}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{M}{\sigma_{M,2}}; \frac{\lambda_0}{\eta} \right) \right] - \frac{\nu^2}{2\eta} - \frac{1}{2\eta\delta} - \frac{\eta}{2} \alpha^2 - \lambda_0 (\varepsilon_0 \sigma_{M,p})^{-q}$$

Note that (E.4) is the dual optimization and hence is a concave problem. We apply the convexity lemma [63, Lemma B.2] to conclude that the objective value in (E.4) also converges to the supremum of function (E.8) over $\lambda_0, \eta \geq 0, \nu$.

Using this observation along with Equation (E.3) and (E.4) we obtain

$$\begin{aligned} & \omega(\alpha, \theta, \varepsilon_0) \\ &= \lim_{n \rightarrow \infty} \omega_s(\mathcal{S}(\alpha, \theta, \varepsilon_0, \boldsymbol{\mu})) \\ &= - \sup_{\lambda_0, \eta \geq 0, \nu} \eta \sqrt{\delta} \mathbb{E} \left[J_q \left(\frac{h}{\eta \sqrt{\delta}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{M}{\sigma_{M,2}}; \frac{\lambda_0}{\eta} \right) \right] - \sqrt{\delta} \left\{ \frac{\nu^2}{2\eta} + \frac{1}{2\eta\delta} + \frac{\eta}{2} \alpha^2 + \lambda_0 (\varepsilon_0 \sigma_{M,p})^{-q} \right\} \\ (E.9) \quad &= \min_{\lambda_0, \eta \geq 0, \nu} \sqrt{\delta} \left\{ \frac{\nu^2}{2\eta} + \frac{1}{2\eta\delta} + \frac{\eta}{2} \alpha^2 + \lambda_0 (\varepsilon_0 \sigma_{M,p})^{-q} \right\} - \eta \sqrt{\delta} \mathbb{E} \left[J_q \left(\frac{h}{\eta \sqrt{\delta}} - \left(\frac{\nu}{\eta} - \theta \right) \frac{M}{\sigma_{M,2}}; \frac{\lambda_0}{\eta} \right) \right], \end{aligned}$$

which completes the proof.

E.3. Proofs that the minimization and maximization primal problems can be restricted to a compact set In this section we demonstrate how the minimization and maximization problems can be restricted to compact sets.

E.3.1. Bounded domains in optimization (B.17) We start with the restriction on $\boldsymbol{\theta}$. Note that one of the claims of Theorem 6.3, part (b), is to show that $\|\tilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2} \rightarrow \alpha_*$ as $n \rightarrow \infty$, in probability, for some α_* by the solution of minimax problem (6.2). We define $\mathcal{S}_\theta = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_{\ell_2} \leq K_\alpha\}$ with $K_\alpha = \alpha_* + \xi$

for a constant $\xi > 0$. We start by the ansatz that $\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}$ and add this ‘artificial constraint’ in the minimax optimization. In addition, by the stationary condition for γ in optimization (B.17) we have

$$\frac{1}{n} \varepsilon \mathbf{u}^T \mathbf{1} = 2\lambda.$$

Therefore

$$\frac{1}{n} \mathbf{u}^T \mathbf{1} = \frac{2\lambda}{\varepsilon} = \frac{2\lambda}{\varepsilon_0 \|\boldsymbol{\mu}\|_{\ell_p}}.$$

Let $\lambda_0 := \frac{\lambda}{\|\boldsymbol{\mu}\|_{\ell_p}}$. Assuming the ansatz that $\lambda_0 = O(1)$, we also use this ‘artificial constraint’ in the minimax optimization.

With these compact constraints in place, we then deploy the CGMT framework to prove Theorem 6.3. This theorem implies that $\|\tilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2} \rightarrow \alpha_*$ as $n \rightarrow \infty$ and so our initial ansatz on the boundedness of $\boldsymbol{\theta}$ is verified. Further, as it can be seen from the proof of Theorem 6.3 (see the line following Equation (B.28)), we have $\lambda_0 = \frac{\lambda}{\|\boldsymbol{\mu}\|_{\ell_p}} \rightarrow \lambda_{0*}$ as $n \rightarrow \infty$, in probability, for some λ_{0*} that is determined by the solution of minimax problem (6.2). This also verifies our ansatz that $\lambda_0 = O(1)$, which in turn implies that $\mathbf{u} \in \mathcal{S}_{\mathbf{u}} = \{\mathbf{u} : 0 \leq u_i, \frac{1}{n} \mathbf{1}^T \mathbf{u} \leq K_{\mathbf{u}}\}$ for some sufficiently large constant $K_{\mathbf{u}} > 0$.

E.3.2. Bounded domains in optimization (B.38) Similar to previous subsection, we start by the ansatz that $\boldsymbol{\theta} \in \mathcal{S}_{\boldsymbol{\theta}}$ where $\mathcal{S}_{\boldsymbol{\theta}} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_{\ell_2} \leq K_{\alpha}\}$ with $K_{\alpha} = \alpha_* + \xi$ for a constant $\xi > 0$, and add this ‘artificial constraint’ in the minimax optimization (B.38). Also by stationarity condition for \mathbf{v} in (B.38) we have $u_i = \ell'(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q})$ and hence $|u_i| \leq K_{\mathbf{u}}$ for some large enough constant $K_{\mathbf{u}} > 0$, using our assumption on the loss function ℓ .

E.4. Proof of Lemma B.1 First note that by Cauchy–Schwarz inequality we have

$$\langle \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{r}, (\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp}) \boldsymbol{\theta}_{\perp} \rangle \geq -\|\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{r}\|_{\ell_2} \left\| (\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp}) \boldsymbol{\theta}_{\perp} \right\|_{\ell_2} = -\alpha \|\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{r}\|_{\ell_2}.$$

To achieve equality, note that similar to (B.50) we have

$$(\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp}) (\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp}) \mathbf{r} = \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{r}.$$

Therefore equality is achieved by choosing $\boldsymbol{\theta} = \lambda \boldsymbol{\Sigma}^{-1/2} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{r}$ with $\lambda = \frac{\alpha}{\|\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{r}\|_{\ell_2}}$.

E.5. Proof of Lemma B.2 By definition of the conjugate function we have

$$\begin{aligned}\tilde{\ell}(\mathbf{v}, \mathbf{w}) &= \sup_{\boldsymbol{\theta}} \mathbf{w}^T \boldsymbol{\theta} - \ell(\mathbf{v}, \boldsymbol{\theta}) \\ &= \sup_{\boldsymbol{\theta}} \mathbf{w}^T \boldsymbol{\theta} - \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q})\end{aligned}$$

Now assume $\boldsymbol{\theta} = \gamma \mathbf{u}$ with $\|\mathbf{u}\|_{\ell_q} = 1$. We thus have,

$$\begin{aligned}\tilde{\ell}(\mathbf{v}, \mathbf{w}) &= \sup_{\mathbf{u}: \|\mathbf{u}\|_{\ell_q}=1, \gamma} \gamma \mathbf{w}^T \mathbf{u} - \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \gamma) \\ &= \sup_{\gamma \geq 0} \gamma \left(\sup_{\mathbf{u}: \|\mathbf{u}\|_{\ell_q}=1} \mathbf{w}^T \mathbf{u} \right) - \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \gamma) \\ &= \sup_{\gamma \geq 0} \gamma \|\mathbf{w}\|_{\ell_p} - \frac{1}{n} \sum_{i=1}^n \ell(v_i - \varepsilon \gamma) .\end{aligned}$$

E.6. Proof of Lemma B.3 By definition

$$\begin{aligned}f^*(\mathbf{u}) &:= \sup_{\tilde{\mathbf{w}}} \langle \mathbf{u}, \tilde{\mathbf{w}} \rangle - f(\tilde{\mathbf{w}}) \\ &= \sup_{\tilde{\mathbf{w}}} \langle \mathbf{u}, \tilde{\mathbf{w}} \rangle + \langle \tilde{\mathbf{w}}, \boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} \rangle \frac{\theta \tau_h}{\alpha} - \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \left\| \boldsymbol{\Sigma}^{1/2} \tilde{\mathbf{w}} \right\|_{\ell_p} \\ &= \sup_{\tilde{\mathbf{w}}} \left\langle \mathbf{u} + \boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} \frac{\theta \tau_h}{\alpha}, \tilde{\mathbf{w}} \right\rangle - \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \left\| \boldsymbol{\Sigma}^{1/2} \tilde{\mathbf{w}} \right\|_{\ell_p} \\ &= \sup_{\tilde{\mathbf{w}}} \left\langle \boldsymbol{\Sigma}^{-1/2} \mathbf{u} + \tilde{\boldsymbol{\mu}} \frac{\theta \tau_h}{\alpha}, \boldsymbol{\Sigma}^{1/2} \tilde{\mathbf{w}} \right\rangle - \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \left\| \boldsymbol{\Sigma}^{1/2} \tilde{\mathbf{w}} \right\|_{\ell_p}\end{aligned}$$

By Hölder's inequality,

$$\left\langle \boldsymbol{\Sigma}^{-1/2} \mathbf{u} + \tilde{\boldsymbol{\mu}} \frac{\theta \tau_h}{\alpha}, \boldsymbol{\Sigma}^{1/2} \tilde{\mathbf{w}} \right\rangle \leq \left\| \boldsymbol{\Sigma}^{-1/2} \mathbf{u} + \tilde{\boldsymbol{\mu}} \frac{\theta \tau_h}{\alpha} \right\|_{\ell_q} \left\| \boldsymbol{\Sigma}^{1/2} \tilde{\mathbf{w}} \right\|_{\ell_p}$$

Therefore, if $\mathbf{u} \in S$ then the supremum is achieved by choosing $\tilde{\mathbf{w}} = 0$. If $\mathbf{u} \notin S$, by scaling $\tilde{\mathbf{w}}$ the supremum would be $+\infty$.

E.7. Proof of Lemma B.4 Fix arbitrary \mathbf{u} . By definition,

$$\begin{aligned}
 \mathbf{P}_{\mathcal{B}}(\mathbf{u}) &:= \arg \min_{\mathbf{z} \in \mathcal{B}} \|\mathbf{u} - \mathbf{z}\|_{\ell_2} \\
 &= \arg \min_{\mathbf{z} \in \mathcal{B}} \|\mathbf{P}_{\boldsymbol{\mu}}^{\perp}(\mathbf{u} - \mathbf{z})\|_{\ell_2}^2 + \|\mathbf{P}_{\boldsymbol{\mu}}(\mathbf{u} - \mathbf{z})\|_{\ell_2}^2 \\
 &\stackrel{(a)}{=} \arg \min_{\mathbf{z} \in \mathcal{B}} \|\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{u} - \mathbf{z}\|_{\ell_2}^2 + \|\mathbf{P}_{\boldsymbol{\mu}} \mathbf{u}\|_{\ell_2}^2 \\
 &= \arg \min_{\mathbf{z} \in \mathcal{B}} \|\mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{u} - \mathbf{z}\|_{\ell_2}^2 \\
 &= \mathbf{P}_{\mathcal{B}} \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{u},
 \end{aligned}
 \tag{E.10}$$

where step (a) follows from that fact that $\mathbf{z} \in \mathcal{B}$ and hence $\mathbf{z} = \mathbf{P}_{\boldsymbol{\mu}}^{\perp} \mathbf{z}$.

E.8. Proof of Lemma B.5 By definition of the set \mathcal{B} , the value of $\|\mathbf{P}_{\mathcal{B}}(\mathbf{h}) - \mathbf{h}\|_{\ell_2}^2$ is given by the optimal objective value of the following optimization:

$$\begin{aligned}
 &\text{minimize}_{\mathbf{z}} \quad \|\mathbf{z} - \mathbf{h}\|_{\ell_2} \\
 &\text{subject to} \quad \left\| \boldsymbol{\Sigma}^{-1/2} \mathbf{z} - \frac{\tau_h \theta}{\alpha} \tilde{\boldsymbol{\mu}} \right\|_{\ell_q} \leq \frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}}, \quad \tilde{\boldsymbol{\mu}}^T \mathbf{z} = 0.
 \end{aligned}
 \tag{E.11}$$

By the change of variable $\mathbf{u} := \boldsymbol{\Sigma}^{-1/2} \mathbf{z} - \frac{\tau_h \theta}{\alpha} \tilde{\boldsymbol{\mu}}$ and forming the Lagrangian, the optimal value of (E.11) is equal to the optimal value of the following problem:

$$\sup_{\lambda \geq 0, \nu} \min_{\mathbf{u}} \frac{1}{2} \left\| \boldsymbol{\Sigma}^{1/2} \left(\mathbf{u} + \frac{\tau_h \theta}{\alpha} \tilde{\boldsymbol{\mu}} \right) - \mathbf{h} \right\|_{\ell_2}^2 + \lambda \left(\|\mathbf{u}\|_{\ell_q}^q - \left(\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \right)^q \right) + \nu \tilde{\boldsymbol{\mu}}^T \boldsymbol{\Sigma}^{1/2} \left(\mathbf{u} + \frac{\tau_h \theta}{\alpha} \tilde{\boldsymbol{\mu}} \right).
 \tag{E.12}$$

Rearranging the terms we get the next alternative representation

$$\sup_{\lambda \geq 0, \nu} \min_{\mathbf{u}} \frac{1}{2} \left\| \boldsymbol{\Sigma}^{1/2} \left(\mathbf{u} + \frac{\tau_h \theta}{\alpha} \tilde{\boldsymbol{\mu}} \right) - \mathbf{h} + \nu \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_{\ell_2}} \right\|_{\ell_2}^2 + \lambda \left(\|\mathbf{u}\|_{\ell_q}^q - \left(\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \right)^q \right) + \nu \tilde{\boldsymbol{\mu}}^T \mathbf{h} - \frac{\nu^2}{2},
 \tag{E.13}$$

Now adopting the notation $\|\mathbf{v}\|_{\boldsymbol{\Sigma}}^2 := \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}$ and invoking the assumption $\boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\mu}} = a \tilde{\boldsymbol{\mu}}$, we rewrite the optimization as follows:

$$\sup_{\lambda \geq 0, \nu} \min_{\mathbf{u}} \frac{1}{2} \left\| \mathbf{u} + \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\boldsymbol{\mu}} - \boldsymbol{\Sigma}^{-1/2} \mathbf{h} \right\|_{\boldsymbol{\Sigma}}^2 + \lambda \left(\|\mathbf{u}\|_{\ell_q}^q - \left(\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \right)^q \right) + \nu \tilde{\boldsymbol{\mu}}^T \mathbf{h} - \frac{\nu^2}{2},
 \tag{E.14}$$

Rearranging the terms further we obtain

$$\begin{aligned}
& \sup_{\lambda \geq 0, \nu} \min_{\mathbf{u}} \left[\frac{1}{2} \left\| \mathbf{u} + \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\boldsymbol{\mu}} - \boldsymbol{\Sigma}^{-1/2} \mathbf{h} \right\|_{\boldsymbol{\Sigma}}^2 + \lambda \|\mathbf{u}\|_{\ell_q}^q \right] - \lambda \left(\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \right)^q + \nu \tilde{\boldsymbol{\mu}}^T \mathbf{h} - \frac{\nu^2}{2} \\
& \text{(E.15)} \\
& = \sup_{\lambda \geq 0, \nu} e_{q, \boldsymbol{\Sigma}} \left(\boldsymbol{\Sigma}^{-1/2} \mathbf{h} - \left(\frac{\tau_h \theta}{\alpha} + \frac{\nu}{a} \right) \tilde{\boldsymbol{\mu}}; \lambda \right) - \lambda \left(\frac{\tau_h}{\alpha} \frac{\gamma_0}{\|\boldsymbol{\mu}\|_{\ell_p}} \right)^q + \nu \tilde{\boldsymbol{\mu}}^T \mathbf{h} - \frac{\nu^2}{2}.
\end{aligned}$$

This concludes the proof.

E.9. Proof of Lemma B.6 We recall the definition of weighted Moreau envelope

$$\text{(E.16)} \quad e_{2, \boldsymbol{\Sigma}}(\mathbf{x}; \lambda) = \min_{\mathbf{v}} \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_{\boldsymbol{\Sigma}}^2 + \lambda \|\mathbf{v}\|_{\ell_2}^2.$$

Setting derivative to zero we get

$$-\boldsymbol{\Sigma}(\mathbf{x} - \mathbf{v}^*) + 2\lambda \mathbf{v}^* = 0,$$

which implies that $\mathbf{v}_* = (\boldsymbol{\Sigma} + 2\lambda \mathbf{I})^{-1} \boldsymbol{\Sigma} \mathbf{x}$. Now consider a singular value decomposition $\boldsymbol{\Sigma} = \mathbf{U} \mathbf{S} \mathbf{U}^T$. Then, $\mathbf{v}_* = \mathbf{U}(\mathbf{S} + 2\lambda \mathbf{I})^{-1} \mathbf{S} \mathbf{U}^T \mathbf{x}$. Substituting for \mathbf{v}_* in (E.16) we obtain

$$\begin{aligned}
e_{2, \boldsymbol{\Sigma}}(\mathbf{x}; \lambda) &= 2\lambda^2 \left\| \mathbf{U}(\mathbf{S} + 2\lambda \mathbf{I})^{-1} \mathbf{S}^{1/2} \mathbf{U}^T \mathbf{x} \right\|_{\ell_2}^2 + \lambda \left\| \mathbf{U}(\mathbf{S} + 2\lambda \mathbf{I})^{-1} \mathbf{S} \mathbf{U}^T \mathbf{x} \right\|_{\ell_2}^2 \\
&= \lambda \mathbf{x}^T \mathbf{U}^T (\mathbf{S} + 2\lambda \mathbf{I})^{-1} \mathbf{S} \mathbf{U} \mathbf{x} \\
&= \lambda \left\| \mathbf{U}(\mathbf{S} + 2\lambda \mathbf{I})^{-1/2} \mathbf{S}^{1/2} \mathbf{U} \mathbf{x} \right\|_{\ell_2}^2 \\
&= \lambda \left\| (\boldsymbol{\Sigma} + 2\lambda \mathbf{I})^{-1/2} \boldsymbol{\Sigma}^{1/2} \mathbf{x} \right\|_{\ell_2}^2
\end{aligned}$$

which yields the desired result.

APPENDIX F: PROPOSITION 3.2 (AN EXTENDED STATEMENT)

Consider the adversarial training loss

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell \left(y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \right),$$

where the loss $\ell(t)$ can be expressed as $\ell(t) = e^{-f(q)}$ obeying the following technical assumptions:

- $f : \mathbb{R} \rightarrow \mathbb{R}$ is C^2 -smooth.
- $f'(q) > 0$ for all $q \in \mathbb{R}$.
- There exists $b_f \geq 0$ such that $qf'(q)$ is non-decreasing for $q \in (b_f, \infty)$ and $qf'(q) \rightarrow \infty$ as $q \rightarrow \infty$.
- Let $g : [f(b_f), \infty) \rightarrow [b_f, \infty)$ be the inverse function of f on the domain $[b_f, \infty)$. There exists $p \geq 0$ such that for all $x > f(b_f)$, $y > b_f$,

$$\left| \frac{g''(x)}{g'(x)} \right| \leq \frac{p}{x}, \quad \left| \frac{f''(y)}{f'(y)} \right| \leq \frac{p}{y}.$$

(It can be verified that the above assumptions are satisfied by exponential loss and logistic loss.) Then, the gradient descent iterates

$$\boldsymbol{\theta}_{\tau+1} = \boldsymbol{\theta}_{\tau} - \mu \nabla \mathcal{L}(\boldsymbol{\theta}_{\tau})$$

with a sufficiently small step size μ obey

$$(F.1) \quad \lim_{t \rightarrow \infty} \left\| \frac{\boldsymbol{\theta}_t}{\|\boldsymbol{\theta}_t\|_{\ell_2}} - \frac{\tilde{\boldsymbol{\theta}}^\varepsilon}{\|\tilde{\boldsymbol{\theta}}^\varepsilon\|_{\ell_2}} \right\|_{\ell_2} = 0,$$

where $\tilde{\boldsymbol{\theta}}^\varepsilon$ is the solution to the following max-margin problem

$$(F.2) \quad \begin{aligned} \tilde{\boldsymbol{\theta}}^\varepsilon = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \quad & \|\boldsymbol{\theta}\|_{\ell_2}^2 \\ \text{subject to} \quad & y_i \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - \varepsilon \|\boldsymbol{\theta}\|_{\ell_q} \geq 1. \end{aligned}$$

UNIVERSITY OF SOUTHERN CALIFORNIA
DATA SCIENCE AND OPERATIONS DEPARTMENT
MARSHALL SCHOOL OF BUSINESS
LOS ANGELES, CA 90089
E-MAIL: ajavanma@usc.edu

UNIVERSITY OF SOUTHERN CALIFORNIA
DEPARTMENT OF ELECTRICAL AND COMPUTER ENGINEERING
VITERBI SCHOOL OF ENGINEERING
LOS ANGELES, CA 90089
E-MAIL: soltanol@usc.edu