THE CURSE OF OVERPARAMETRIZATION IN ADVERSARIAL TRAINING: PRECISE ANALYSIS OF ROBUST GENERALIZATION FOR RANDOM FEATURES REGRESSION

BY HAMED HASSANI^{1,a}, ADEL JAVANMARD^{2,b}

Successful deep learning models often involve training neural network architectures that contain more parameters than the number of training samples. Such overparametrized models have recently been extensively studied, and the virtues of overparametrization have been established from both the statistical perspective, via the double-descent phenomenon, and the computational perspective via the structural properties of the optimization landscape. Despite this success, it is also well known that these models are highly vulnerable to small adversarial perturbations in their inputs. Even when adversarially trained, their performance on perturbed inputs (robust generalization) is considerably worse than their best attainable performance on benign inputs (standard generalization). It is thus imperative to understand how overparametrization fundamentally affects robustness.

In this paper, we will provide a precise characterization of the role of overparametrization on robustness by focusing on random features regression models (two-layer neural networks with random first layer weights). We consider a regime where the sample size, the input dimension and the number of parameters grow proportionally, and derive an asymptotically exact formula for the robust generalization error when the model is adversarially trained. Our developed theory reveals the nontrivial effect of overparametrization on robustness and indicates that high overparametrization can hurt robust generalization.

1. Introduction The success of deep learning models is often reliant on training highly complex neural networks whose number of parameters is much larger than the number of data points. Even though the large complexity of such models allows for perfect interpolation of the data, they often achieve low generalization error. This behavior has resulted in a growing body of work aiming to analyze such so-called overparametrized models.

Recent work has demonstrated the virtues of overparametrization from statistical and optimization-based perspectives. From the statistical viewpoint, it is now well-documented that many overparametrized models exhibit a 'double-descent' property [5, 4, 61]: As the model complexity increases, the generalization error first follows the traditional U-shaped curve until a specific point, after which the error decreases, and attains a global minimum in the overparametrized regime. In fact, the minimum generalization error often appears to be at infinite complexity – the more overparametrized is the model, the smaller is the error. It is often argued that the good generalization behavior of highly overparametrized models is due to the inductive bias of gradient-based algorithms which helps with selecting models that generalize well –see e.g, [3, 37, 79, 35]. From the optimization viewpoint, training deep neural networks in general involves optimizing highly non-convex functions, but it has been conjectured that

¹Department of Electrical and Systems Engineering, University of Pennsylvania, ^ahassani@seas.upenn.edu

²Data Sciences and Operations Department, University of Southern California, ^bajayanma@usc.edu

MSC2020 subject classifications: Primary 62E20, 62F12; secondary 62F35.

Keywords and phrases: adversarial training, random features models, precise high-dimensional asymptotics, Gaussian equivalence property.

highly overparametrized models are easy to optimize despite non-convexity. Instances of this observation has been formally proved, e.g. in [77, 44, 41, 3, 64]. The high-level intuition here is that in the highly overparametrized regimes, a model that perfectly interpolates the training data (and so is a global minimizer of the empirical risk) is found in the neighborhood of most initializations.

Despite the remarkable success of deep neural networks, and the crucial role of over-parametrization in both the generalization and the tractability aspects, these models are known to be highly vulnerable to perturbations in the input [6, 84]. With an unguarded training approach, these models show unsatisfactory *robust generalization error* in the presence of "small" worst-case perturbations to their inputs, a.k.a *asdversarial examples*. This suggests that learning algorithms, even those with excellent performance on test data, may not be learning the true underlying concepts that determine the response; although they work well on naturally occurring data, adversarial examples have low probability in the data distribution and expose fundamental blind spots in the learning algorithms.

This observation stimulated significant effort to improve robustness using a wide variety of *adversarial training* methods which often involve augmenting the training loss so as to become more robust to input perturbations (see e.g. [32, 47, 43, 57, 93, 97, 13]). However, there is still a large gap between the robust generalization error and the (standard) generalization error in adversarially-trained models. In summary, while modern overparametrized machine learning models perform very well on benign inputs, they still remain fragile to perturbations in the input. These findings raise a fundamental question:

How does overparametrization affect robustness to perturbations in the input?

A few recent papers have begun to answer the above question in specific settings with rather conflicting messages:

- [46] and [23] have studied high-dimensional *linear* models and showed that the robust generalization error of adversarially-trained models becomes *worse* as the models become more overparametrized. It should be noted that for linear models, even in the case where there is no adversary, the best generalization error is attained when the model is underparametrized [37].
- Another line of work provably shows that in order to *interpolate* the training data smoothly, while being robust, overparametrization is *necessary* [10, 9]. However, we note that, in order to train robust models, it may not be beneficial to interpolate the training data as robustness is measured with worst-case performance over all the points in a neighborhood around the input data. Indeed, [22] study the tradeoffs between memorization (of training data) and robustness of two-layer neural networks and established a lower-bound on the non-robustness of the model (via the Sobolev-seminorm of the model) as an increasing function of the amount of memorization.

We will provide a more detailed discussion of these points and other related works in Section 3. Despite such interesting recent progress, a comprehensive understanding on how overparametrization precisely affects robustness remains largely mysterious.

In this paper, we focus on random features regression models that are adversarially trained using robust empirical risk minimization and provide a "precise characterization" of the robust generalization. Our analysis is carried out in a high-dimensional regime where the size of the training data n, the number of parameters N, and the dimension of the data d grow proportional to each other, i.e. $N/d \rightarrow \psi_1$ and $n/d \rightarrow \psi_2$. We further assume that the perturbations are bounded in terms of ℓ_2 norm by a value $\varepsilon > 0$. Our developed theory allows us to precisely characterize the effect of overparamterization on model robustness. One of the

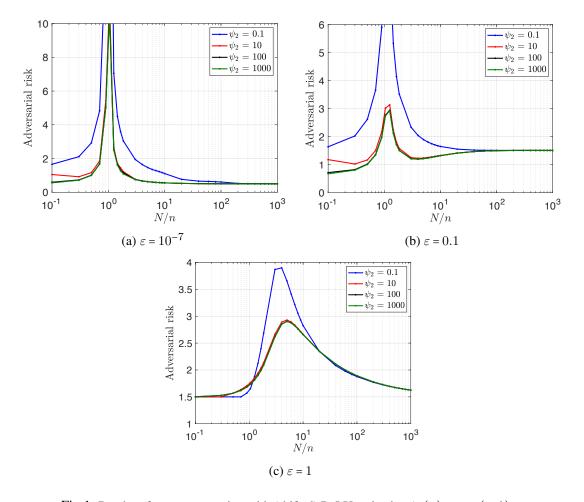


Fig 1: Random features regression with (shifted) ReLU activation $(\sigma(x) = \max(x,0) - 1/\sqrt{2\pi})$. Data (x_i,y_i) is generated with d-dimensional normal covariates x_i and $y_i = \beta^T x_i + \xi_i$, where the noise variables $\xi_i \sim \mathcal{N}(0,0.5)$ and $\|\beta\|_{\ell_2} = 1$. Perturbations are allowed within an Euclidean ball of radius ε , and the models are adversarially trained. We plot the robust generalization error (using Theorem 4.2) versus the amount of overparametrization N/n, where N is the number of parameters and n is the number of training data points. The plots are obtained for different values of ε and $\psi_2 = n/d$.

main consequences of our analysis is that, in general, higher overparametrization leads to a *worse* robust generalization error for the adversarially-trained models. Figure 1 depicts how the robust generalization error varies with respect to the amount of overparametrization N/n. The left figure corresponds to the case where there is no adversary (i.e. $\varepsilon \approx 0$). In this case, the robust generalization coincides with the (standard) generalization error, and is minimized at infinite overparametrization. However, as seen in the other two figures (for $\varepsilon > 0$), overparametrization is in general hurting robustness. This is clearly seen in Figure 1(c) and (b) (for $\psi_2 \ge 10$) where the minimum robust error is attained when the model is underparametrized. We refer to Figure 4 for an extended version of Figure 1 with more choices of ε and signal-to-noise ratios.

We proceed by providing an informal overview of our results and their implications in Section 2. Related works are discussed in Section 3. The main result of the paper, which characterizes the robust generalization error for the random features model is explained in Section 4. The architecture of the proof of the main result is sketched in Section 6. Our analysis develops a set of techniques that are of independent interest: (i) We derive an asymptotic

closed form for adversarial examples in trained random features models; (ii) While features are highly non-Gaussian in random features models, we prove a Gaussian equivalence property which relates robust generalization in these models to that of linear models with Gaussian features under the same correlation structure; (iii) Our analysis of the equivalent Gaussian model relies on the Convex Gaussian Min-max Theorem, which is a generalized and tight version of Gordon's Gaussian comparison inequalities.

2. Results and discussion: An informal overview

Problem setting. Consider a supervised learning scenario where we are given i.i.d data $\{(x_i, y_i)\}_{i \le n}$ generated according to the following distribution:

$$(2.1) y_i = \langle \boldsymbol{x}_i, \boldsymbol{\beta} \rangle + \xi_i, \quad \text{with} \quad \boldsymbol{x}_i \sim_{iid} \mathsf{N}(0, \boldsymbol{I}_d), \quad \xi_i \sim \mathsf{N}(0, \tau^2).$$

The (linear) dependence between (x_i, y_i) is unknown and the goal is to fit a model to this data which can be then used to predict labels for the unlabeled examples at test time.

We consider modeling the relation between label y and feature vector x using the class of random features (RF) model, which can be described as

(2.2)
$$\mathcal{F}_{RF}(\boldsymbol{W}) = \left\{ f(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{W}) = \sum_{\ell=1}^{N} \theta_{\ell} \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle) : \boldsymbol{\theta} = (\theta_{1}, \dots, \theta_{N}) \in \mathbb{R}^{N} \right\},$$

where $\boldsymbol{\theta}$ is the parameter vector to be learned and $\boldsymbol{W} \in \mathbb{R}^{N \times d}$ is a fixed matrix whose rows \boldsymbol{w}_{ℓ} are chosen randomly and independently of data. For simplicity we assume the normalization $\|\boldsymbol{w}_{\ell}\|_{\ell_2} = 1$. Namely, the vectors \boldsymbol{w}_{ℓ} are chosen uniformly at random from the unit sphere, $\boldsymbol{w}_{\ell} \sim \mathrm{Unif}(\mathbb{S}^{d-1})$, which implies that $\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{j} \rangle$ is of order one. In addition, $\sigma : \mathbb{R} \mapsto \mathbb{R}$ is a nonlinear activation function.

Note that in random features model training is only done on θ and not on W. In other words, the random features model can be perceived as a two-layer neural network with the weights of the first layer chosen randomly and independently from data, while the weights in the second layer are learned during the training phase. The random features model was introduced by [68] 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 [15, 14, 42, 51]. The random features models are arguably the simplest analytically tractable models that capture all the features of the double descent phenomenon without assuming ad hoc misspecification structures [63]. In particular, they allow to disentangle the number of parameters from the covariates dimension and hence isolate the effects of overparametrization from the effects of the ambient dimension.

To quantify robust generalization, we consider an adversarial framework where at the test time, the feature vector x is corrupted by additive perturbation, chosen adversarially, from the Euclidean ball of radius ε . We measure the robust generalization via the *adversarial risk* measure which is the expected test error of the model on perturbed test input. We train a random features model, using a widely used adversarial training approach, which is based on the robust empirical risk minimizer (robust-ERM estimator) [58, 90]:

(2.3)
$$\widehat{\boldsymbol{\theta}}^{\varepsilon} = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^{N}} \max_{\|\boldsymbol{\delta}_{i}\|_{\ell_{2}} \le \varepsilon} \frac{1}{2n} \sum_{i=1}^{n} (y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}))^{2}.$$

where δ_i is the norm-bounded adversarial perturbation on sample covariates x_i and ε is the "perceived" adversary's power used in the training process.

Results and discussion. We study the asymptotic setting, where $N, n, d \to \infty$ with $N/d \to \psi_1$ and $n/d \to \psi_2$ for some positive constants ψ_1, ψ_2 . We derive the precise characterization of

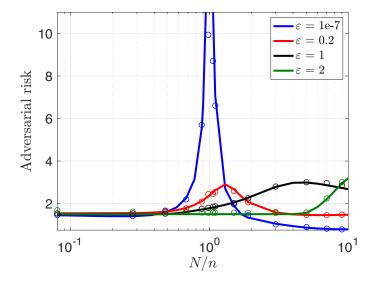


Fig 2: Adversarial risk versus overparametrization $\psi_1/\psi_2=N/n$ for different values of adversary's power ε_0 . Solid curves are theoretical predictions and dots are results obtained based on gradient descent on the robust ERM objective. Each dot represents the average of 100 trials. The data is generated according to model (2.1), with d=100, n=300, $\tau^2=0.5$, and $\boldsymbol{\beta}\in\mathbb{R}^d$ obtained by drawing a vector with i.i.d N(0,1) entries and then normalizing it to have $\|\boldsymbol{\beta}\|_{\ell_2}=1$.

the adversarial risk of the robust-ERM estimator, as an explicit function of the dimension parameters ψ_1, ψ_2 , the noise level τ^2 , and the adversarial power ε . We refer to Theorem 4.2 for the specific formulae.

Let us now discuss the behavior of the robust generalization curve under different settings. We consider the data model (2.1) and the random features regression with shifted ReLU activation:

$$\sigma(x) = \max(x,0) - \frac{1}{\sqrt{2\pi}}.$$

The reason behind the intercept term is that since the response variable is zero mean, we consider fitting a model using zero mean features. Note that $\langle w_\ell, \boldsymbol{x} \rangle \sim \mathsf{N}(0,1)$ and for $G \sim \mathsf{N}(0,1)$, we have $\mathbb{E}[\sigma(G)] = \mathbb{E}[G\mathbb{I}(G>0) - 1/\sqrt{2\pi}] = 0$.

We start by Figure 2 which shows our theoretical curve versus the overparametrization ratio $\psi_1/\psi_2=N/n$ along with the corresponding empirical results. The solid lines depict theoretical predictions with the dots representing the empirical performance of gradient descent in learning the robust ERM for data model (2.1), with $d=100,\,n=300,\,\tau^2=0.5$. In addition, $\boldsymbol{\beta}\in\mathbb{R}^d$ is generated by first drawing a d-dimensional vector with i.i.d standard normal entries and then normalizing it to have unit ℓ_2 norm. Each dot represents the average of 20 trials. As we see, even for moderate covariate dimensions (d), our theoretical curve is at excellent match with the empirical results. We note that when $\varepsilon \to 0$ (we did not set $\varepsilon=0$ exactly for numerical stability), we are in non-adversarial regime and the robust generalization error reduces to the usual test error (blue curve). In this case, we observe the double-descent phenomena and recover the theoretical prediction of [61]. As ε grows the robust generalization curve starts behaving differently. For ε large enough ($\varepsilon=1,2$ in the figure), we see that overparametrization hurts robust generalization.

For a more complete picture, in Figure 3 we consider similar setting with more choices of ε and noise variance τ^2 , and also a larger range of overparametrization $\psi_2/\psi_1 = N/n$, as we fix

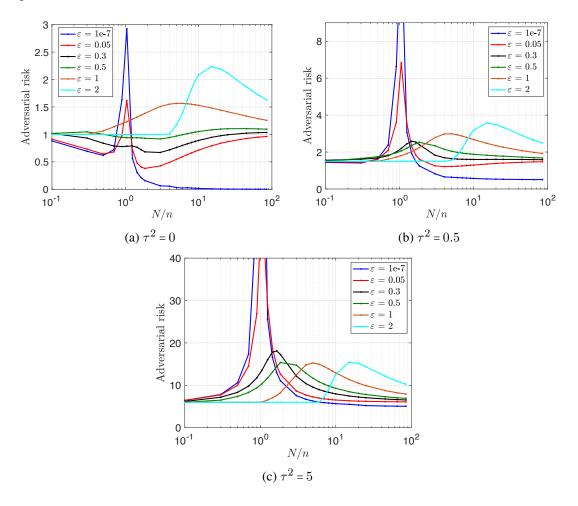


Fig 3: Theoretical prediction curves for adversarial risk of robust ERM as a function of overparametrization $\psi_1/\psi_2 = N/n$ for different values of adversary's power ε and noise variance τ^2 , with the data model (2.1). Here we fix $\|\boldsymbol{\beta}\|_{\ell_2} = 1$ and $\psi_2 = 3$.

 ψ_2 = 3. When $N/n \to 0$, we essentially have the risk of the zero estimator, which is $\|\beta\|_{\ell_2}^2 + \tau^2$. Several intriguing observations can be made from these plots:

- In the noiseless case (Figure 3a) and for $\varepsilon \le 0.5$, the global minimizer of the adversarial risk is at a finite overparametrization (N/n > 1), after which the risk becomes increasing as a function of N/n (higher overparametrization hurst robust generalization). Similar behavior is observed for $\tau^2 = 0.5$ and $\varepsilon \le 0.05$.
- In all three plots (corresponding to different SNR levels), when ε is large enough ($\varepsilon \ge 1$), the risk first goes up as overparametrization increases and after reaching its peak starts going down, but it remains above $1 + \tau^2$ which is the risk at the highly underparametrized regime $(N/n \to 0)$. Therefore, somewhat surprisingly, robust ERM estimator has larger adversarial risk compared to the trivial zero estimator, for all the range of overparametrization.
- The peak of the adversarial risk occurs in the overparametrized regime; for the non-adversarial case $\varepsilon = 0$, it occurs at the interpolation threshold N/n = 1 and for $\varepsilon > 0$ it occurs at N/n > 1. The location of the peak and the value of risk at the peak vary with ε . As ε grows, the peak shifts to the right and occurs at a higher overparametrization ratio.

In Figure 4 we depict our theoretical prediction curves for the adversarial risk of the robust ERM estimator as a function of the overparametrization ratio $\psi_1/\psi_2 = N/n$ for different values

of $\psi_2 = n/d$. The right panel corresponds to $\varepsilon = 1$ (strong adversary) and as we see for different values of τ^2 and ψ_2 , the adversarial risk is first an increasing function of overparametrization ratio, until it reaches its peak (in the overparametrized regime, N/n > 1) and then becomes decreasing. But it never falls below its initial value at $N/n \approx 0$. The left panel corresponds to $\varepsilon = 0.1$ (weak adversary) and as we see for large ψ_2 , overparametrization clearly has a negative effect on robust generalization. For example, in Figure 4a, for $\psi_2 = 100, 1000$ the risk is an increasing function of ψ_1/ψ_2 over the entire range. Also in Figure 4c, for $\psi_2 \ge 10$ the global minimum of the adversarial risk is achieved in the underparametrized regime (N/n < 1).

3. Related Work Several recent works have focused on the robustness of overparametrized models. On the one hand, [9] shows that in order *interpolate* the training data smoothly, the Lipschitz parameter of the resulting model should be at least of order $\sqrt{\frac{nd}{N}}$. This applies to data distributions that satisfy a property called isometry–e.g. when the data covariates \boldsymbol{x}_i are distributed on the unit sphere. For such data distributions, worst-case perturbations are meaningful only if their ℓ_2 norm is upper-bounded by $\frac{\epsilon}{\sqrt{d}}$. Otherwise, if the size of the perturbation can be allowed to be much larger than $O(\frac{1}{\sqrt{d}})$, it can be shown that the robust generalization error approaches one for any model–see [75, 29, 59, 60]. Putting the above two results together, we can conclude that, in order to interpolate smoothly, while guaranteeing robustness to norm-bounded perturbations, it is necessary that the ratio N/n is bounded away from zero. This is indeed the regime studied in our paper. However, in this regime, it is not clear why interpolation to training data is beneficial for obtaining robust models. In fact, to obtain robust models, one may have to trade off the performance on the original data points (i.e. interpolation to training data) with the performance on the points in a ball around each data point (i.e. extrapolation to adversarial examples). In other words, we may have underparametrized models that do not fit the training data perfectly, but have a small Lipschitz constant. Indeed, this can be implied from the main messages of our paper.

On the other hand, the works in [46] and [23] have studied the performance of highdimensional *linear* models and showed that the robust generalization error of adversariallytrained models becomes worse as the models become more overparametrized. In particular, [23] provably shows that avoiding interpolation (and using underparametrized models) improves the robust generalization error in both linear regression and classification—which leads to the first theoretical result on robust overfitting. There are a few reasons on why we might prefer to study non-linear models (such as the random features model) compared to linear models [61]: First of all, for linear models, we know that the best (standard) generalization error is attained when the model is highly underparametrized. Second, the number of parameters in a linear model is tied to the covariates dimension d and hence the effects of overparametrization cannot be isolated from the effects of the ambient dimensions. Third, a hypothesis put forward in [32] is that the origins and ubiquity of adversarial examples is due to the (approximately) linear behavior of a model over large regions of the input space. Shallow linear models are not able to become constant near training points while also assigning different outputs to different training points. However, the setting of random features is significantly different since this class can express any function to an arbitrary degree of accuracy so long as it has enough number of random features [69].

In another recent work [94], an extensive study on the robustness of wide neural networks with respect to norm-bounded perturbations is provided. By defining and analyzing a new metric, called perturbation stability, it is shown that while the (standard) generalization error is improved on wider models, the perturbation stability often worsens, leading to a potential decrease in the overall model robustness. These empirical findings are aligned with the messages of our paper.

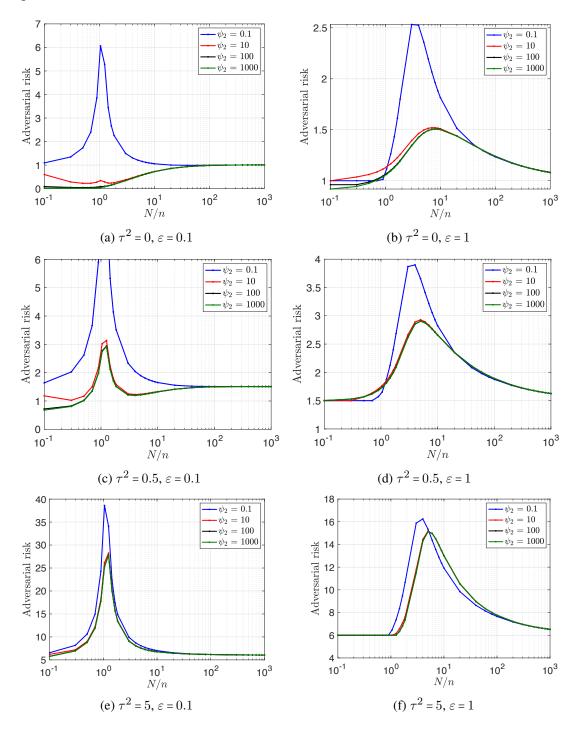


Fig 4: Theoretical prediction curves for adversarial risk of robust ERM as a function of overparametrization $\psi_1/\psi_2 = N/n$ for different values of $\psi_2 = n/d$. Each plot corresponds to a specific value of adversary's power ε and noise variance τ^2 .

A somehow different line of work [96] studies the sample complexity of the robust interpolation problem where the goal is to interpolate (noisy) training data by a Lipschitz function, under generic covariate distribution (beyond isoperimetry distributions). This work measures the (non)robustness of a model by its Lipschitz constant and similar to [22] establish a lower

bound on Lipschitnzenss which is increasing with respect to the overfitting level. This result can be rephrased as an adverse effect of overparametrization (thorough overfitting) on robustness. While these work study the effect of overparametrization on robustness via memorization/interpolation, we will take a direct approach to study the effect of overparametrization on 'adversarially trained' models.

Several works have shown a non-trivial tradeoff between the robust generalization error and the standard generalization error for parametric models [90, 83, 67, 97, 46, 21]. It has also been shown that using more data can improve this tradeoff [11, 62, 67, 74, 17, 95, 65, 34, 70]. Again, these findings are aligned with the messages of our paper as more data can mean less overparametrization.

This paper provides, for the first time, an analysis for the adversarially-trained random features model in the high-dimensional regime. For linear models, this analysis has been carried out in [46] for the regression setting, and later on in [45, 85] for the classification setting. A key ingredient of the analysis in these papers, as well as our paper, is a powerful extension of a classical Gaussian process inequality [33], known as the Convex Gaussian Minimax Theorem, developed in [88] and further extended in [87, 16]. Another key ingredient of our analysis is the Gaussian Equivalence Property for the random features model which was proposed and studied in [63] for maximum-margin linear classifiers in the overparametrized regime, as well as [37, 1, 63, 61, 28, 19, 39, 31, 30] for the linear Gaussian model under other settings. In particular, a part of our analysis, that establishes equivalence with the so-called noisy linear model in the adversarial setting, is heavily based on the machinery which was elegantly developed in [39] for the random features model. This machinery is itself based on the Lindeberg principle [55] and the leave-one-out technique developed in [27, 1].

We conclude this section by a broad comparison between adversarial setting and the literature of robust statistics.

Comparison with robust statistics. This area traditionally considers a setting where perturbations are made to the *training data*; a small fraction of data samples are grossly corrupted and the goal is to find estimators that are robust against outliers (via measures like influence function, breakdown point, and change of variance, etc). In the adversarial training paradigm, one considers the so-called test-time adversarial setting, in which the training data is uncorrupted (say $(x_i, y_i) \sim \mathbb{P}$ for some distribution \mathbb{P}). However, the adversary can perturb *each test data*. In other words, the test data (x, y) is drawn from \mathbb{P} and then x is perturbed by the adversary $(x \to \tilde{x})$ with $\|x - \tilde{x}\|_{\ell_2} \le \varepsilon$. The goal of adversarial training is to develop a model that can still predict the response y from the perturbed feature \tilde{x} . With this view, adversarially robust models are basically those that have good generalization and are also smooth enough so that they do not change much on small neighborhoods (of radius at most ε).

That said, another line of work (see e.g. [48]) 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. Related to this setting, there is also a line of work based on the Median of Means approach, see e.g, [40, 18]), which concerns a data poisoning/ data contamination adversarial setting. In data poisoning, the adversary can pick a (small) fraction of the *training data* and alter it in a way that it hurts the training process, and ultimately the generalization performance. However, in this paper we consider a different type of adversarial act which has to do with adversarial perturbation (in a small ball) of the input data point at the *test time*.

4. Main results Recall the data distribution given in (2.1). Given n i.i.d pairs (x_i, y_i) drawn from this distribution, we fit a random features model, defined as the function class (2.2), with the shifted ReLU activation:

(4.1)
$$\sigma(x) = \max(x,0) - \frac{1}{\sqrt{2\pi}}.$$

We consider sequences of parameters (N, n, d) that diverge proportionally to each other and sometimes, we index such sequences by d, with N = N(d) and n = n(d) functions of d.

Assumption 1 (Asymptotic setting.)

(a) Defining $\psi_{1,d} = N/d$ and $\psi_{2,d} = n/d$, we assume that the following limits exist:

$$\lim_{d\to\infty}\psi_{1,d}=\psi_1,\quad \lim_{d\to\infty}\psi_{2,d}=\psi_2\,,$$

for some positive finite constants ψ_1 and ψ_2 .

(b) We assume that the ℓ_2 norm of the signal $\boldsymbol{\beta}$ converges, as $d \to \infty$. For the sake of normalization and without loss of generality, we assume $\lim_{d\to\infty} \|\boldsymbol{\beta}\|_{\ell_2} = 1$.

Recall that in the data model (2.1), $x_i \sim_{iid} N(0, I_d)$ and so its distribution is rotation-invariant. Likewise, in the random features model (2.2), the rows w_ℓ are chosen uniformly at random from the unit sphere, and so has a rotation-invariant distribution. In our adversarial setting, we also focus on norm-bounded perturbations which is again a rotation-invariant constraint. Using these properties, it is easy to see that the adversarial risk will be invariant if we rotate the model β in (2.1) and hence only depends on $\|\beta\|_{\ell_2}$. This justifies Assumption 1(b) made above.

To study robust generalization of the estimated models, we consider an adversarial framework with norm bounded perturbations. This can be formulated as a game between the learner and the adversary. Given access to a training dataset consisting of n i.i.d. pairs (x_i, y_i) generated from (2.1), the learner chooses a model θ from the class of random features model $\mathcal{F}_{RF}(W)$ (2.2). Adversarial perturbations occur at the test time. After observing the learner's model, for every test point x, the adversary perturbs it to $x + \delta$ where δ is chosen from the Euclidean ball of radius ε . Note that the choice of δ can in general depend on x and the learner's model. The robust generalization of the learner's model is quantified via a measure called *adversarial risk*, which is the expected prediction loss of the model on an adversarially corrupted test data point according to the described attack model.

Definition 4.1 (Adversarial risk.) For a predictive model f and a loss of choice $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}_{\geq 0}$, the adversarial risk of model f is defined as:

$$\mathsf{AR}(f) \coloneqq \mathbb{E}\Big[\max_{\|\pmb{\delta}\|_{\ell_2} \leq \varepsilon_{\text{test}}} \ell(f(\pmb{x} + \pmb{\delta}), y)\Big],$$

where the expectation is with respect to randomness of (x, y).

In particular, for a random features model $\boldsymbol{\theta} = (\theta_1, \dots, \theta_N)^\mathsf{T}$ from $\mathcal{F}_{RF}(\boldsymbol{W})$, defined in (2.2), and with the choice of squared loss, the adversarial risk of $\boldsymbol{\theta}$ becomes:

(4.2)
$$\mathsf{AR}(\boldsymbol{\theta}) \coloneqq \mathbb{E}\left[\max_{\|\boldsymbol{\delta}\|_{\ell_2} \le \varepsilon_{\text{test}}} \left(y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}(\boldsymbol{x} + \boldsymbol{\delta}))\right)^2\right].$$

Norm bounded adversarial attack models are widely used in the literature, motivated by a plethora of safety-critical applications in machine learning, computer vision, natural language processing, medical imaging, and robotics. A popular approach to adversarial training is by considering the following robust empirical risk minimization (robust-ERM) problem [58, 90]:

(4.3)
$$\widehat{\boldsymbol{\theta}}^{\varepsilon} = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^{N}} \max_{\|\boldsymbol{\delta}_{i}\|_{\ell_{2}} \leq \varepsilon} \frac{1}{2n} \sum_{i=1}^{n} (y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}))^{2}.$$

Here, $\varepsilon_{\rm test}$ is a measure of the adversary's power and ε is the "perceived" adversary's power used by the algorithm. Our theory allows for ε to be different from $\varepsilon_{\rm test}$; cf Theorem 4.2. In our numerical experiments in Section 2 we consider $\varepsilon = \varepsilon_{\rm test}$ to focus on other relevant quantities, namely ψ_1 , ψ_2 on adversarial risk.

Note that the above objective is the empirical surrogate of the adversarial risk (4.2), where the expectation is replaced by the empirical average over the training samples. This minimax approach can also be viewed as an implicit smoothing that tries to fit the same response y to all the covariate vectors in the ε -neighborhood of x simultaneously.

Our main result in this paper is a precise characterization of the adversarial risk of the robust ERM model (4.3) under the asymptotic regime described in Assumption 1. Before stating our result, we introduce another piece of notation.

For $\psi_1 \in (0, \infty)$, we define function $S(\cdot; \psi_1) : \mathbb{R}_{<0} \to \mathbb{R}_{<0}$:

(4.4)
$$S(z;\psi_1) = \frac{1 - \psi_1 - z - \sqrt{(1 - \psi_1 - z)^2 - 4\psi_1 z}}{-2\psi_1 z}.$$

One may recognize that $S(z; \psi_1)$ is the Stieltjes transform of the Marchenko-Pastur distribution. We refer to Lemma F.2 (Appendix F) for more details.

We are now ready to state our main result.

Theorem 4.2 Let n i.i.d pairs (x_i, y_i) be drawn from the data model (2.1) and let $\widehat{\theta}^{\varepsilon}$ be the robust ERM fit (4.3) to this data using the class of random features models $\mathcal{F}_{RF}(W)$, given by (2.2) with the shifted ReLU activation. Consider the asymptotic regime, described in Assumption 1. With function $S(\cdot; \psi_1)$ given by (4.4), define

$$\sigma^2 = \tau^2 + 1 - \psi_1 \left(1 + \left(1 - \frac{2}{\pi} \right) S\left(\frac{2}{\pi} - 1; \psi_1 \right) \right).$$

(a) For $\varepsilon > 0$, the following convex-concave minimax scalar optimization has a unique solution $(\alpha_*, \tau_{q*}, \beta_*, \gamma_*, \tau_{q*})$:

(4.5)
$$\max_{0 \leq \beta, \gamma, \tau_q} \min_{0 \leq \alpha, \tau_g} \mathcal{R}(\alpha, \tau_g, \beta, \gamma, \tau_q),$$

where

$$\mathcal{R}(\alpha, \tau_{g}, \beta, \gamma, \tau_{q}) \coloneqq \frac{\tau_{q}}{2\alpha} (\tau^{2} + 1 - \sigma^{2}) - \frac{\alpha \tau_{q}}{2} + \frac{\beta \tau_{g}}{2} \psi_{2} + \frac{\beta}{2(\tau_{g} + \beta)} (\sigma^{2} + \alpha^{2})$$

$$+ \mathbf{1}_{\left\{\frac{\gamma(\tau_{g} + \beta)}{\varepsilon \beta \sqrt{\alpha^{2} + \sigma^{2}}} > \sqrt{\frac{2}{\pi}}\right\}} \frac{\beta^{2} (\alpha^{2} + \sigma^{2})}{2\tau_{g} (\tau_{g} + \beta)} \left(\operatorname{erf} \left(\frac{\nu^{*}}{\sqrt{2}} \right) - \frac{\gamma(\tau_{g} + \beta)}{\varepsilon \beta \sqrt{\alpha^{2} + \sigma^{2}}} \nu^{*} \right)$$

$$- \frac{\alpha}{\tau_{q}} \sup_{0 \leq \lambda < 1} \left[\frac{\lambda \psi_{1}}{2} \left\{ \frac{\tau_{q}^{2}}{\alpha^{2}} + \beta^{2} + \left(\frac{\tau_{q}^{2}}{\alpha^{2}} \left(1 - \frac{2}{\pi} \lambda \right) + \frac{2}{\pi} (1 - \lambda) \beta^{2} \right) S\left(\frac{2}{\pi} \lambda - 1; \psi_{1} \right) \right\} - \frac{\lambda}{2(1 - \lambda)} \gamma^{2} \right].$$

Here, ν^* is the unique solution to

$$\frac{\gamma(\tau_g + \beta)}{\varepsilon \beta \sqrt{\alpha^2 + \sigma^2}} - \frac{\beta}{\tau_g} \nu - \nu \cdot \operatorname{erf}\left(\frac{\nu}{\sqrt{2}}\right) - \sqrt{\frac{2}{\pi}} e^{-\frac{\nu^2}{2}} = 0.$$

(b) The adversarial risk of the robust ERM $\widehat{\theta}^{\varepsilon}$ converges in probability

(4.6)
$$\mathsf{AR}(\widehat{\boldsymbol{\theta}}^{\varepsilon}) \overset{\mathcal{P}}{\to} \left[1 + \left(\frac{\varepsilon_{\text{test}} \beta_{\star} \nu_{\star}}{\varepsilon \tau_{g\star}} \right)^{2} + 2 \sqrt{\frac{2}{\pi}} \frac{\varepsilon_{\text{test}} \beta_{\star} \nu_{\star}}{\varepsilon \tau_{g\star}} \right] (\alpha_{\star}^{2} + \sigma^{2}).$$

Here, the probabilistic statement is with respect to the randomness in both the training data $\{(x_i, y_i)\}_{i=1}^n$ and the random features W.

We note that the robust ERM estimator is a random and rather complicated high-dimensional function of the training data. However, in the asymptotic regime where $N, n, d \rightarrow \infty$ at the same order, the adversarial risk of the robust ERM concentrates and the above theorem provides an exact characterization of its limit as a deterministic formula. The derived formula is based on a five dimensional convex-concave mini-max optimization problem and its optimal solution can be easily derived via a simple low-dimensional gradient descent rather quickly and accurately. Alternatively, one can form a system of equations by writing the KKT stationary conditions corresponding to (4.5). The adversarial risk prediction can then be written in terms of the fixed point of this system of deterministic equations.

Let us re-emphasize the contribution of theorem 4.2. Albeit its involved form, it describes the behavior of a *high-dimensional random* problem in terms of a *deterministic* optimization with a handful number of scalar variables. This theme of result is similar to state evolution equation for approximate message passing algorithms [24], density evolution for LDPC codes [71, Chapter 4], and characterizing the trajectory of SGD for training neural networks in terms of partial differential equations [78, 44].

Remark 4.1 (Solving Optimization (4.5)) We find the solution to this optimization by solving for the first-order optimality conditions (stationary equations). We set the (sub)gradient with respect to the variables to zero since it is non-smooth, which results in a system of nonlinear equations with seven variables, namely $\alpha, \tau_g, \beta, \gamma, \tau_q, \lambda, \nu^*$. In the numerical experiments, we use fsolve command in Matlab to solve this system of equations, which is based on the trust-region algorithm.

- **5. Discussion** By virtue of Theorem 4.2 we characterize the adversarial risk of the robust ERM in term of the solution of the deterministic optimization problem (4.5). Given that it does not admit a closed-form solution in general, and is rather involved, in this section we discuss some of the applications of this theorem including optimal choice of ε during training, trend of adversarial risk with respect to different quantities, and implications for non-adversarial setting.
- 5.1. Optimal ε for robust ERM estimator: An interesting application of our theory is to derive the optimal ε (perceived adversary's perturbation level) in the robust ERM, while fixing the adversary's (actual) perturbation level on test inputs to $\varepsilon_{\text{test}}$. The optimal ε here refers to the value which minimized the adversarial risk. An intriguing observation is that the optimal ε is different than $\varepsilon_{\text{test}}$ in general, and depends on ψ_1, ψ_2 in a non-trivial way (There is no universal solution, which underscores the significance of possessing a precise theory that comprehends the impact of adversarial training, which constitutes the principal objective of the present work.)

In Figure 5a, we fix $\psi_2 = 3$, $\tau = \sqrt{0.5}$, $\varepsilon_{\text{test}} = 0.3$, and plot the adversarial risk of $\widehat{\theta}^{\varepsilon}$ as we vary ε for different values of ψ_1 . As we see the optimal value of ε (resulting in minimum risk) changes with ψ_1 , it is in general different from the test adversary's perturbation $\varepsilon_{\text{test}}$. In addition, the optimal ε increases with ψ_1 . In Figure 5b, we plot similar curves, fixing $\varepsilon_{\text{test}} = 0$. In this case, the adversarial risk reduces to the notion of standard risk defined as

(5.1)
$$SR(\boldsymbol{\theta}) := \mathbb{E}\left[(y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}))^{2} \right].$$

As we see form the plots, even though $\varepsilon_{\text{test}} = 0$, adversarial training can help as minimum risk is achieved at positive ε . The reason is that adversarial training acts as a regularization (It becomes clearer after derivation (6.11), where adversarial training aims to find solution with small $\|J\theta\|_{\ell_2}$.) In particular, we observe that

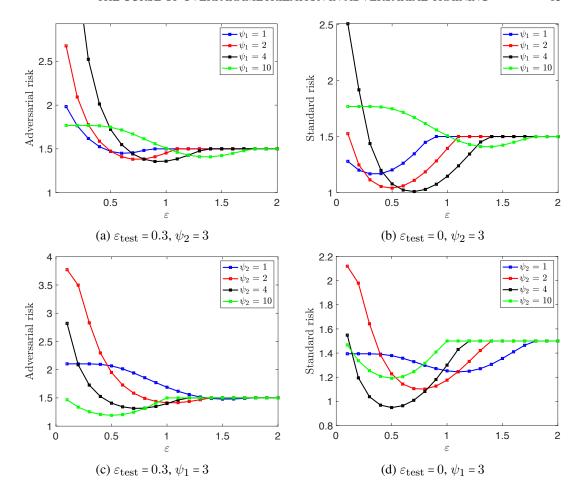


Fig 5: Behavior of adversarial/standard risk as we vary ε , the "perceived" adversary's power used in the adversarial training. In (a), (b), ψ_2 = 3 is fixed and in (c), (d), we fix ψ_1 = 3. Also, (b), (d) correspond to $\varepsilon_{\text{test}}$ = 0, and so there is no perturbation at the test time. In these cases, adversarial risk reduces to the standard risk. In these experiments, we set τ^2 = 0.5, $\|\beta\|_{\ell_2}$ = 1.

- The optimal ε is always greater than or equal to $\varepsilon_{\rm test}$, the true test perturbation magnitude. This 'additional' regularization helps with minimizing the adversarial risk.
- At higher overparametrization measured by $\psi_1/\psi_2 = N/n$ the benefit of this regularization becomes stronger. This is also evident from Figure 5c, where with fixed ψ_2 . As we increase ψ_1 , the optimal ε also increases. Likewise, in Figure 5d, with fixed ψ_1 , the optimal ε increases as ψ_2 decreases.

In summary, our theory allows to understand when adversarial training is beneficial and what is the optimal value of ε to use in training (depending on $\psi_1, \psi_2, \varepsilon_{\text{test}}$ and τ .)

- 5.2. Dependence on ε , ψ_1 , ψ_2 : We discern the following trends in the analytical curves for adversarial risk which are derived based on Theorem 4.2 (Equation (4.6)).
- From Figure 5, fixing $\varepsilon_{\rm test}$, ψ_1 and ψ_2 , the adversarial risk is first decreasing in ε until it gets to its minimum, after which it becomes increasing in ε indicating that the regularization effect from adversarial training is larger than it should be, and then eventually it levels

at $\sigma^2 + \|\boldsymbol{\theta}_0\|^2$ (the risk of model $\boldsymbol{\theta} = \mathbf{0}$, which is the ERM when the regularization is very strong).

- For fixed ψ_2 , if SNR is large enough and ε is small enough, we observe a double descent behavior in the adversarial risk (see Figure 3a, $\varepsilon = 1e 7$ and $\varepsilon = 0.05$ and Figure 4c).
- Fixing $\varepsilon_{\text{test}}$ and ψ_1 , the adversarial risk becomes decreasing in ψ_2 after some point. In the next subsection, we characterize the adversarial risk in non-adversarial setting, by which we see this trend for $\psi_2 > \psi_1$.
- 5.3. Non-adversarial training: An important special case of the result is when $\varepsilon = \varepsilon_{test} = 0$. In words, there is no adversarial-training and also there is no adversary during the test time. Theorem 4.2 allows us to characterize the standard risk (generalization error) of the ERM estimator.

We focus on the underparameterized regime (n > N) or equivalently $\psi_2 > \psi_1$, since otherwise at $\varepsilon = 0$ the problem is underdetermined. In this case the objective \mathcal{R} in (4.5) significantly simplifies and we can indeed obtain a closed-form expression for the risk.

Proposition 5.1 Let n i.i.d pairs (x_i, y_i) be drawn from the data model (2.1) and let $\widehat{\theta}$ be the ERM (4.3) fit to this data using the class of random features models $\mathcal{F}_{RF}(W)$, given by

$$\widehat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^N} \frac{1}{2n} \sum_{i=1}^n (y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i))^2,$$

with $\sigma(\cdot)$ the shifted ReLU activation. Consider the asymptotic regime, described in Assumption 1. With function $S(\cdot; \psi_1)$ given by (4.4), define

$$\sigma^2 = \tau^2 + 1 - \psi_1 \left(1 + \left(1 - \frac{2}{\pi} \right) S\left(\frac{2}{\pi} - 1; \psi_1 \right) \right).$$

Then, the standard risk of the ERM $\widehat{\theta}$ converges in probability to

(5.2)
$$\operatorname{SR}(\widehat{\boldsymbol{\theta}}) \stackrel{\mathcal{P}}{\to} \sigma^2 \left(\frac{\psi_2}{\psi_2 - \psi_1}\right).$$

We refer to Section E.4 for the proof of this proposition.

We next use the result of Proposition 5.1 to discuss the role of ψ_1 and ψ_2 on the risk:

- Recall that σ only depends on τ and ψ_1 . Fixing ψ_1 the dependence of risk on ψ_2 is of form $\psi_2/(\psi_2-\psi_1)$. This is decreasing in $\psi_2=n/d$. For example if d,N are fixed, and we increase the sample size n, the risk goes down which is expected.
- Dependence on ψ_1 is more involved as the term σ^2 also depends on ψ_1 through the Stieltjes transform S. In Figure 6 we plot the risk versus ψ_1 for different values of ψ_2 . As we see up to some threshold, it is decreasing in ψ_1 but after that it becomes increasing. This is expected because for example fixing n, d (and so ψ_2), as we increase N (and so ψ_1), first the risk goes down because the model becomes richer to capture the data generative model, but after some point it has a reverse effect, because we need to estimate larger number of parameters N, from fixed sample size n, while this excess model complexity is not needed. As we see in the plots, this threshold is increasing with ψ_2 .
- It is also worth comparing the standard risk of random features model with that of linear models. For $\psi_2 \ge 1$, using the result of [37, Proposition 2], the risk of ridgeless least squares is given by $\tau^2 \psi_2/(\psi_2 1)$. This is similar to our characterization (5.2), where the noise variance τ^2 is replaced with the effective noise variance σ^2 , and ψ_2 is replaced by $\psi_2/\psi_1 = n/N$. (Note that the number of parameters to be learnt in the linear model is d, while in the random features model is N. So the sample size to parameter size ratio in the linear regression is ψ_2 , while for the random features model it is ψ_2/ψ_1 .)

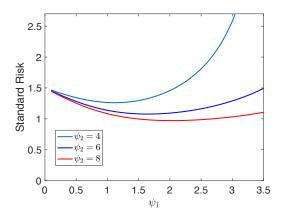


Fig 6: Behavior of the standard risk of the ERM estimator $\hat{\theta}$ versus ψ_1 for different values of ψ_2 . As observed, the risk is first decreasing in ψ_1 , up to some threshold depending on ψ_2 , after which it becomes increasing. This threshold is increasing with ψ_2 .

6. Architecture of the proof This section introduces the key steps underlying the proof of our main result, Theorem 4.2. Our analysis is intricate and consists of a host of novel ideas which could be of separate interest. Here we discuss the major steps, along with an overview of the techniques and intermediate results.

Define the loss $\mathcal{L}(\theta)$ given by

(6.1)
$$\mathcal{L}(\boldsymbol{\theta}) \coloneqq \max_{\|\boldsymbol{\delta}_i\|_{\ell_2} \le \varepsilon} \frac{1}{2n} \sum_{i=1}^n (y_i - \boldsymbol{\theta}^\mathsf{T} \sigma (\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i))^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \Omega \boldsymbol{\theta},$$

where $\Omega := I + \frac{\sqrt{\log(d)}}{d} \mathbf{1} \mathbf{1}^{\mathsf{T}}$. Here, $\mathbf{1} = (1, 1, \dots, 1) \in \mathbb{R}^N$ and $\zeta > 0$ is an arbitrary small but fixed constant. By definition, the robust ERM estimator (4.3) is the minimizer of $\mathcal{L}(\boldsymbol{\theta})$ for $\zeta = 0$.

The regularization $\frac{\zeta}{2}(\|\boldsymbol{\theta}\|_{\ell_2}^2 + \frac{\sqrt{\log(d)}}{d}\langle \mathbf{1}, \boldsymbol{\theta}\rangle^2)$ in the loss $\mathcal{L}(\boldsymbol{\theta})$ is added for technical reasons. In our analysis, we let $\zeta \to 0$, after letting $d \to \infty$ to characterize the adversarial risk of the robust ERM $\widehat{\boldsymbol{\theta}}^{\varepsilon}$. We refer to Section A in the supplementary for the justification of this step.

Before we outline the main steps of the proof, we note that since the rows of the matrix $W \in \mathbb{R}^{N \times d}$ are generated i.i.d. according to $w_{\ell} \sim \mathrm{Unif}(\mathbb{S}^{d-1})$, then the matrix norm of W is bounded with high probability and the rows of W are almost orthogonal. More precisely, we define the event

(6.2)
$$\mathcal{E}_{\boldsymbol{W}} := \left\{ \|\boldsymbol{W}\| \le \sqrt{\psi_{1,d}} + C, \ \left|\boldsymbol{w}_{\ell}^{\mathsf{T}} \boldsymbol{w}_{k}\right| \le \log(d) / \sqrt{d} \quad \forall \ell \ne k \right\},$$

for a large enough constant C (note that $N/d := \psi_{1,d}$ – see Assumption 1). Using well-known results on the norm of random matrices (see e.g. [91, Theorem 5.39]) as well as Hoeffding's inequality we have $\mathbb{P}(\mathcal{E}_{W}) \ge 1 - c \exp(-\log^2(d)/c)$ for some constant c > 0. In the following, our statements are proven conditioned on the event \mathcal{E}_{W} which holds with high probability.

Notation. We need to define a few pieces of notation which will be used in the following. We use $O_d(\cdot)$, $o_d(\cdot)$ to denote the standard big-O and little-o notation, where we stress the asymptotic variable d. Likewise, we use $O_{d,\mathbb{P}}$ and $o_{d,\mathbb{P}}$ to indicate asymptotic behavior in probability. Specifically, $f(d) = O_{d,\mathbb{P}}(g(d))$ if for any $\varepsilon > 0$, there exists $C_{\varepsilon} > 0$ and large enough d_{ε} such that $\mathbb{P}(|f(d)/g(d)| > C_{\varepsilon}) \le \varepsilon$, for all $d \ge d_{\varepsilon}$. Similarly, $f(d) = o_{d,\mathbb{P}}(g(d))$ if f(d)/g(d) converges to zero in probability. We write $f(d) \approx g(d)$ as $d \to \infty$, when f(d) = f(d)/g(d) is f(d)/g(d).

 $g(d) \to 0$, in probability. Note that we consider the asymptotic regime where n,d,N grow at the same scale, $(\lim N/d \to \psi_1 \text{ and } \lim n/d \to \psi_2 \text{ for some positive constants } \psi_1 \text{ and } \psi_2)$, the expression $d \to \infty$ implies that $n,N \to \infty$, as well.

For a matrix A, we denote by $\|A\|$ its operator norm, $\|A\|_F = (\sum_{ij} A_{ij}^2)^{1/2}$ the Frobenius norm of A. For an integer n, we use the shorthand $[n] = \{1, \ldots, n\}$.

Finally, the indicator function is denoted by $\mathbb{I}(\cdot)$ – i.e., $\mathbb{I}(A)$ = 1 only if the event A holds true, and otherwise $\mathbb{I}(A)$ = 0.

6.1. An asymptotically-exact closed form for adversarial examples We start by simplifying the loss $\mathcal{L}(\theta)$. If the activation function σ was linear, then finding the worst-case perturbations δ_i (maximizers in the definition of loss $\mathcal{L}(\theta)$) amounts to a trust-region subproblem that can be solved in closed form—see [46]. A major challenge here is the nonlinearity of σ . In the first step we use the specific form of the activation to derive an asymptotically equivalent but simpler form of $\mathcal{L}(\theta)$.

Note that $\sigma(z) = \max(z,0) - 1/\sqrt{2\pi}$ is linear on the positive z and constant on the negative z. Also by the constraint on perturbation we have $\|\langle \boldsymbol{w}, \boldsymbol{\delta} \rangle\| \leq \|\boldsymbol{w}\|_{\ell_2} \|\boldsymbol{\delta}\|_{\ell_2} \leq \varepsilon$. Therefore, if $\|\langle \boldsymbol{w}, \boldsymbol{x} \rangle\| \geq \varepsilon$, then $\langle \boldsymbol{w}, \boldsymbol{x} \rangle$ and the perturbed form $\langle \boldsymbol{w}, \boldsymbol{x} + \boldsymbol{\delta} \rangle$ share the same sign. In this case, the worst case $\boldsymbol{\delta}$ can be solved exactly. One can also use the randomness in \boldsymbol{x} to bound the number of rows of \boldsymbol{W} for which $|\langle \boldsymbol{w}, \boldsymbol{x} \rangle| < \gamma$, where γ is some small constant, and show that the contribution of these terms in the loss is asymptotically negligible. This argument is formalized in the next proposition. All the proofs of the statements in this section are relegated to Appendix B.

Proposition 6.1 Assume (x,y) generated according to (2.1). Further define

(6.3)
$$C_{\boldsymbol{\theta}} \coloneqq \{ \boldsymbol{\theta} \in \mathbb{R}^N : \| \boldsymbol{\theta} \|_{\ell_{\infty}} \le C_0 \sqrt{\log(d)/d}, \| \boldsymbol{\theta} \|_{\ell_2} \le C_0 \},$$

for an arbitrary but fixed constant $C_0 > 0$. Then, we have

(6.4)

$$\max_{\|\boldsymbol{\delta}\|_{\ell_2} \le \varepsilon} \left| y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}(\boldsymbol{x} + \boldsymbol{\delta})) \right| = \left| y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}) \right| + \varepsilon \left\| \boldsymbol{W}^\mathsf{T} \mathrm{diag} \left(\mathbb{I}(\boldsymbol{W} \boldsymbol{x} > 0) \right) \boldsymbol{\theta} \right\|_{\ell_2} + O_{d, \mathbb{P}} \left(\frac{\log(d)}{d^{1/6}} \right),$$

uniformly over $\theta \in C_{\theta}$. Here, the probability bound is with respect to the randomness in x. I.e. W is fixed and event \mathcal{E}_{W} in (6.2) is assumed to hold.

To be able to use the result of above proposition, we show that the minimizer of $\mathcal{L}(\theta)$ falls in \mathcal{C}_{θ} defined in (6.3).

Proposition 6.2 Assume (x,y) is generated according to (2.1), and recall that the rows of $\mathbf{W} \in \mathbb{R}^{N \times d}$ are drawn i.i.d. from $\mathrm{Unif}(\mathbb{S}^{d-1})$. Let $\widehat{\boldsymbol{\theta}} = \arg\min \mathcal{L}(\boldsymbol{\theta})$. We have $\widehat{\boldsymbol{\theta}} \in \mathcal{C}_{\boldsymbol{\theta}}$, with probability at least $1 - 4e^{-cn}$ for some absolute constant c > 0.

Motivated by Proposition 6.1 we define loss $\overset{\circ}{\mathcal{L}}(\theta)$ as follows:

$$(6.5) \quad \overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}) \coloneqq \frac{1}{2n} \sum_{i=1}^{n} \left(|y_i - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{W}^{\mathsf{T}} \operatorname{diag} \left(\mathbb{I}(\boldsymbol{W} \boldsymbol{x}_i > 0) \right) \boldsymbol{\theta} \|_{\ell_2} \right)^2 + \frac{\zeta}{2} \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{\theta}.$$

By using Proposition 6.1, we can prove the following.

Proposition 6.3 *Under the setting of Proposition 6.1 we have*

(6.6)
$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \left| \mathcal{L}(\boldsymbol{\theta}) - \overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}) \right| = o_{d,\mathbb{P}}(1).$$

6.2. Concentration of the adversarial effects As can be observed from the loss $\mathcal{L}(\boldsymbol{\theta})$, the effect of adversarial perturbation in reflected via the terms $\eta_i \coloneqq \| \boldsymbol{W}^\mathsf{T} \mathrm{diag}(\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)) \boldsymbol{\theta} \|_{\ell_2}$. In the next proposition, we show that apart from a negligible fraction of data points $i \in [n]$, the perturbation terms $\eta_i(\boldsymbol{\theta})^2$ concentrate around their expectation. All the proofs of the statements in this section are relegated to Appendix C.

Proposition 6.4 Let $\eta_i(\boldsymbol{\theta}) \coloneqq \| \boldsymbol{W}^\mathsf{T} \mathrm{diag}(\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)) \boldsymbol{\theta} \|_{\ell_2}$ and $\nu_i(\boldsymbol{\theta}, \gamma) \coloneqq \mathbb{I}(|\eta_i(\boldsymbol{\theta})|^2 - \mathbb{E}[\eta_i(\boldsymbol{\theta})|^2]| > \gamma)$. Under the setting of Proposition 6.1 and for any sequence γ_d such that $1/\gamma_d = e^{o(\sqrt{\log(d)})}$, we have

(6.7)
$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \frac{1}{n} \sum_{i=1}^{n} \nu_i(\boldsymbol{\theta}; \gamma_d) = O_{d, \mathbb{P}}(\log(d\gamma_d^2)^{-0.5}).$$

Corollary 6.5 By choosing sequence $\gamma_d = 1/\log(d)$ we obtain

(6.8)
$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \frac{1}{n} \sum_{i=1}^{n} \nu_i(\boldsymbol{\theta}; \frac{1}{\log(d)}) = o_{d, \mathbb{P}}(1).$$

By Corollary 6.5, other than at most an $o_d(1)$ fraction of data points $i \in [n]$, the terms $\eta_i(\theta)^2$ concentrate, in the sense $|\eta_i(\theta)^2 - \mathbb{E}[\eta_i(\theta)^2]| \le 1/\log(d)$, uniformly over $\theta \in \mathcal{C}_{\theta}$. This suggests that in the loss function we can replace the terms $\eta_i(\theta)$ by $\sqrt{\mathbb{E}[\eta_i(\theta)^2]}$. This observation will be formally stated in the next lemma. Before proceeding to it, let us compute the expectation of terms $\eta_i(\theta)^2$. We write

(6.9)

$$\mathbb{E}[\eta_i(\boldsymbol{\theta})^2] = \mathbb{E}\left[\|\boldsymbol{W}^\mathsf{T} \operatorname{diag}(\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0))\boldsymbol{\theta}\|_{\ell_2}^2\right] = \boldsymbol{\theta}^\mathsf{T} \,\mathbb{E}\left[(\boldsymbol{W}\boldsymbol{W}^\mathsf{T}) \odot \left(\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)^\mathsf{T}\right)\right]\boldsymbol{\theta},$$

where the expectation is with respect to x_i and W is fixed. Hence, this can be written as $\mathbb{E}[\eta_i(\theta)^2] = \|J\theta\|_{\ell_2}^2$ with

$$\boldsymbol{J} \coloneqq ((\boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}) \odot \mathbb{E}[\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)^{\mathsf{T}}])^{1/2}$$
.

Note that the J is well-defined since the matrix under the square root is positive semidefinite. (This follows from the observation that the expression (6.9) is positive for all θ .)

Since $\|\boldsymbol{w}_{\ell}\|_{\ell_2} = 1$ and $\boldsymbol{x}_i \sim \mathsf{N}(0, \boldsymbol{I}_d)$, we have that $\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_i \rangle$ and $\langle \boldsymbol{w}_k, \boldsymbol{x}_i \rangle$ are jointly Gaussian with

$$\mathbb{E}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle^{2}) = \mathbb{E}(\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle^{2}) = 1, \quad \mathbb{E}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle \langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle) = \langle \boldsymbol{w}_{k}, \boldsymbol{w}_{\ell} \rangle,$$

and by using [15, Table 1] we get

$$\mathbb{E}[\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)^{\mathsf{T}}] = \frac{\pi - \cos^{-1}(\boldsymbol{W}\boldsymbol{W}^{\mathsf{T}})}{2\pi}.$$

Therefore, we obtain the following explicit formulation for J_W :

(6.10)
$$\mathbf{J} = \left((\mathbf{W} \mathbf{W}^{\mathsf{T}}) \odot \left(\frac{\pi - \cos^{-1} (\mathbf{W} \mathbf{W}^{\mathsf{T}})}{2\pi} \right) \right)^{1/2}.$$

Motivated by Corollary 6.5 and the interpretation after it we define the loss function

(6.11)
$$\overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}) \coloneqq \frac{1}{2n} \sum_{i=1}^{n} \left(|y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2} \right)^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta}.$$

By our next lemma, the minimizer of $\overset{\circ}{\mathcal{L}}(\theta)$ converges to the minimizer of the original loss $\mathcal{L}(\theta)$ and therefore we can work with $\overset{\circ}{\mathcal{L}}(\theta)$ for our asymptotic analysis.

Lemma 6.6 We have

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \frac{|\overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta})|}{1 + \min(\mathcal{L}(\boldsymbol{\theta}), \overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}))} = o_{d,\mathbb{P}}(1).$$

Also, by denoting by $\widehat{\boldsymbol{\theta}}^*$ and $\widehat{\boldsymbol{\theta}}$ the minimizers of $\mathcal{L}(\boldsymbol{\theta})$ and $\mathcal{L}(\boldsymbol{\theta})$, we have $\|\widehat{\boldsymbol{\theta}}^* - \widehat{\boldsymbol{\theta}}\|_{\ell_2} \to 0$, in probability.

Motivated by the result of Lemma 6.6, we define a notion of adversarial risk based on the modified loss $\overset{\circ}{\mathcal{L}}(\theta)$. Specifically, we define

(6.12)
$$\overset{\circ\circ}{\mathsf{AR}}(\boldsymbol{\theta}) \coloneqq \mathbb{E}_{\boldsymbol{x},y} \left[\left(|y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x})| + \varepsilon_{\text{test}} \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_2} \right)^2 \right].$$

In the next lemma, we show that $\stackrel{\circ}{\mathsf{AR}}(\theta)$ converges to $\mathsf{AR}(\theta)$ uniformly over \mathcal{C}_{θ} .

Lemma 6.7 Recall the adversarial risk of a model θ , denoted by $AR(\theta)$ and given by (4.2). Let $\stackrel{\circ}{AR}(\theta)$ be defined as (6.12). We then have,

$$\sup_{\boldsymbol{\theta} \in \mathsf{C}_{\boldsymbol{\theta}}} \frac{|\mathsf{AR}(\boldsymbol{\theta}) - \mathsf{AR}(\boldsymbol{\theta})|}{\sqrt{\mathsf{AR}(\boldsymbol{\theta})}} = o_{d,\mathbb{P}}(1).$$

6.3. The Gaussian equivalence property and the noisy linear model In this section, we will show that in order to characterize the robust generalization error of the random features model, we can equivalently consider the so-called Gaussian features model (a.k.a. the noisy linear model). This equivalency is often termed as the Gaussian Equivalence Property (GEP), and has recently been proven in several contexts [63, 39, 28, 19]. We prove this equivalency for adversarially-trained random features models in this section.

We begin with decomposing the nonlinear activation function $\sigma(z)$ as

(6.13)
$$\sigma(z) = \mu_0 + \mu_1 z + \mu_2 \sigma_+(z),$$

where for $G \sim N(0,1)$,

$$\mu_0 := \mathbb{E}[\sigma(G)], \quad \mu_1 = \mathbb{E}[G\sigma(G)], \quad \mu_2 := \sqrt{\mathbb{E}[\sigma^2(G)] - \mu_0^2 - \mu_1^2}.$$

For the case of shifted ReLU activation, defined in (4.1), we have $\mu_0 = 0$, $\mu_1 = \frac{1}{2}$ and $\mu_2 = \sqrt{\frac{1}{4} - \frac{1}{2\pi}}$. Also, $\sigma_{\perp}(z)$ is the nonlinear component of the activation function which is orthogonal to the constant and linear components in the following sense: $\mathbb{E}[\sigma_{\perp}(G)] = 0$ and $\mathbb{E}[G\sigma_{\perp}(G)] = 0$. We can then write the random features $\sigma(\boldsymbol{W}\boldsymbol{x})$ as follows

(6.14)
$$\sigma(\mathbf{W}\mathbf{x}) = \mu_0 \mathbf{1} + \mu_1 \mathbf{W}\mathbf{x} + \mu_2 \sigma_1 (\mathbf{W}\mathbf{x}),$$

Note that the random variables $\sigma(\boldsymbol{w}_i^{\mathsf{T}}\boldsymbol{x})$ have zero mean and unit variance, by construction. Further, $\mathbb{E}_{\boldsymbol{x}}\{(\boldsymbol{w}_i^{\mathsf{T}}\boldsymbol{x})\sigma_{\perp}(\boldsymbol{w}_i^{\mathsf{T}}\boldsymbol{x})\}=0$ since by construction $\mathbb{E}[\sigma_{\perp}(G)G]=0$. This suggests to replace the variables $\sigma_{\perp}(\boldsymbol{w}_i^{\mathsf{T}}\boldsymbol{x})$ by a set of i.i.d standard normal variables and consider the following *noisy linear model*

(6.15)
$$f := \mu_0 \mathbf{1} + \mu_1 \mathbf{W} \mathbf{x} + \mu_2 \mathbf{u},$$

with $f, u \in \mathbb{R}^N$ and $u \sim N(0, I_N)$ is generated independently from x.

Consequently, we define the loss of the noisy linear model as

(6.16)
$$\mathcal{L}_{\mathrm{nl}}(\boldsymbol{\theta}) \coloneqq \frac{1}{2n} \sum_{i=1}^{n} (|y_i - \boldsymbol{\theta}^\mathsf{T} \boldsymbol{f}_i| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta},$$

where f_i are generated i.i.d. according to (6.15). We note that compared to the loss $\mathcal{L}(\theta)$ defined in (6.11), we have only replaced the feature vectors $\sigma(\boldsymbol{W}\boldsymbol{x}_i)$ with the noisy linear features \boldsymbol{f}_i .

Let $\widehat{\boldsymbol{\theta}}^*$ and $\widehat{\boldsymbol{\theta}}_{nl}^*$ respectively denote the minimizers of $\mathcal{L}(\boldsymbol{\theta})$ and $\mathcal{L}_{nl}(\boldsymbol{\theta})$. Roughly speaking, the Gaussian equivalence property (GEP) states that under certain conditions on \boldsymbol{W} and the activation function σ , we have

(6.17)
$$\overset{\circ}{\mathsf{AR}}(\widehat{\boldsymbol{\theta}}^*) \approx \overset{\circ}{\mathsf{AR}}_{\mathrm{nl}}(\widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^*) \quad \text{as } d \to \infty,$$

where $\overset{\circ\circ}{AR}(\cdot)$ is defined by (6.12) and $\overset{\circ\circ}{AR}_{nl}(\cdot)$ is its counterpart defined based on the noisy linear model, as follows:

(6.18)
$$\overset{\circ}{\mathsf{AR}_{\mathrm{nl}}}(\boldsymbol{\theta}) \coloneqq \mathbb{E}_{\boldsymbol{f},y} \left[\left(|y - \boldsymbol{\theta}^\mathsf{T} \boldsymbol{f}| + \varepsilon_{\mathrm{test}} \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_2} \right)^2 \right].$$

Therefore, by virtue of Lemma 6.7 and (6.17), we can henceforth focus on characterizing $\stackrel{\circ}{\mathsf{AR}}_{\mathrm{nl}}(\widehat{\theta}_{\mathtt{nl}}^*)$.

In order to prove (6.17), we first show the asymptotic equality of $\overset{\circ}{AR}_{nl}(\theta)$ and $\overset{\circ}{AR}(\theta)$. All the proofs of the statements in this section are provided in Appendix D.

Proposition 6.8 Consider model (2.1) under the asymptotic setting in Assumption 1 and define the set

$$C'_{\boldsymbol{\theta}} \coloneqq \left\{ \boldsymbol{\theta} : \|\boldsymbol{\theta}\|_{\ell_{\infty}} \le C_0 \sqrt{(\log(d))/d} \text{ and } \|\boldsymbol{\theta}\|_{\ell_2} \le C_0 \text{ and } |\mathbf{1}^{\mathsf{T}}\boldsymbol{\theta}| \le C_0 \sqrt{d/(\log(d))} \right\},$$

for an arbitrary but fixed constant $C_0 > 0$. Let $\overset{\circ\circ}{\mathsf{AR}}(\boldsymbol{\theta})$ and $\overset{\circ\circ}{\mathsf{AR}}_{\mathrm{nl}}(\boldsymbol{\theta})$ be defined by (6.12)-(6.18). Then, for any $\boldsymbol{\theta} \in \mathcal{C}'_{\boldsymbol{\theta}}$ we have

(6.19)
$$\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta}) = \overset{\circ}{\mathsf{AR}}_{\mathrm{nl}}(\boldsymbol{\theta}) + o_d(1),$$

In addition, we have the following characterizations for $\overset{\circ\circ}{\mathsf{AR}}_{\mathrm{nl}}(\boldsymbol{\theta})$:

(6.20)
$$\overset{\circ \circ}{\mathsf{AR}_{\mathrm{nl}}}(\boldsymbol{\theta}) = M(\boldsymbol{\theta})^2 + \varepsilon_{\mathrm{test}}^2 \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_2}^2 + 2\sqrt{\frac{2}{\pi}}\varepsilon_{\mathrm{test}}M(\boldsymbol{\theta}) \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_2} ,$$

with $M(\boldsymbol{\theta})$ given by

(6.21)
$$M(\boldsymbol{\theta})^{2} = \tau^{2} + \left\| \frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_{2}}^{2} + \left(\frac{1}{4} - \frac{1}{2\pi} \right) \|\boldsymbol{\theta}\|_{\ell_{2}}^{2}.$$

Proof of Proposition 6.8, i.e. equation (6.19), follows from a Central limit theorem (CLT) for weakly correlated variables proved in [30]. Specifically, [30] shows that $(\boldsymbol{\theta}^{\mathsf{T}}\sigma(\boldsymbol{W}\boldsymbol{x}), \boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x})$ converges in distribution to $(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{f},\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x})$, where $\boldsymbol{\beta}$ is a fixed vector with bounded norm. In [39], the authors provide an alternative proof of this CLT using Stein's method and the Lindeberg approach. Their analysis assumes that the activation function $\sigma(z)$ is an odd function with bounded first derivatives. In addition, their analysis gives the convergence rate in terms of $\|\boldsymbol{\theta}\|_{\ell_{\infty}}$ (a Berry-Esseen type result).

By the characterizations (6.20), we know that, provided that $\theta \in \mathcal{C}'_{\theta}$, the quantity $\operatorname{AR}_{\operatorname{nl}}(\theta)$ depends on θ through the quantities $M(\theta)$ and $\|J\theta\|_{\ell_2}$. As we show in Lemma F.3, $\|J^2 - K\| \to 0$, in probability with $K = (WW^{\mathsf{T}} + I)/4$. Since by definition, for $\theta \in \mathcal{C}'_{\theta}$ we have $\|\theta\|_{\ell_2} \le C_0$, therefore $\|J\theta\|_{\ell_2}^2 \to (\|W^{\mathsf{T}}\theta\|_{\ell_2}^2 + \|\theta\|_{\ell_2}^2)/4$. So, in order to show the GEP relation of the form (6.17), it suffices to show that the quantities $\|\theta\|_{\ell_2}$, $\|\frac{1}{2}W^{\mathsf{T}}\theta - \beta\|_{\ell_2}$ and $\|W^{\mathsf{T}}\theta\|_{\ell_2}$ evaluated at $\widehat{\theta}^*$ converge to the corresponding quantities evaluated at $\widehat{\theta}^*_{\mathrm{nl}}$, and also $\widehat{\theta}^*$, $\widehat{\theta}^*_{\mathrm{nl}} \in \mathcal{C}'_{\theta}$.

Theorem 6.9 Consider the quantities Φ_A and Φ_B defined as

$$\Phi_A := \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} (|y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2})^2 + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \|\frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta}\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta})^2,$$

$$\Phi_{B} \coloneqq \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} (y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f}_{i} + \varepsilon \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}})^{2} + \lambda \|\boldsymbol{\theta}\|_{\ell_{2}}^{2} + \lambda_{w} \|\frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta}\|_{\ell_{2}}^{2} + \lambda_{s} \frac{\log(d)}{d} (\boldsymbol{1}^{\mathsf{T}} \boldsymbol{\theta})^{2},$$

where $\lambda, \lambda_s, \lambda_w > 0$, and (x_i, y_i) is generated i.i.d. according to (2.1). We further assume that the event $\mathcal{E}_{\mathbf{W}}$ holds. Then, we have

(6.22)
$$\Phi_A \xrightarrow{\mathcal{P}} c \text{ if and only if } \Phi_B \xrightarrow{\mathcal{P}} c,$$

where $\stackrel{\mathcal{P}}{\longrightarrow}$ denotes convergence in probability.

Using this theorem, we can then prove the following proposition.

Proposition 6.10 Recall $\widehat{\theta}^*$ and $\widehat{\theta}_{nl}^*$ given by

$$\widehat{\boldsymbol{\theta}}^* = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^N} \stackrel{\circ}{\mathcal{L}}(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^N} \frac{1}{2n} \sum_{i=1}^n (|y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \Omega \boldsymbol{\theta},$$

(6.23)
$$\widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^* = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^N} \mathcal{L}_{\mathrm{nl}}(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^N} \frac{1}{2n} \sum_{i=1}^n (|y_i - \boldsymbol{\theta}^\mathsf{T} \boldsymbol{f}_i| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta}.$$

Then, under the asymptotic regime of Assumption 1 we have

$$\widehat{\boldsymbol{\theta}}^*, \widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^* \in \mathcal{C}_{\boldsymbol{\theta}}',$$

with probability $1 - o_d(1)$, and

(6.24)
$$M(\widehat{\boldsymbol{\theta}}^*) - M(\widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^*) \xrightarrow{\mathcal{P}} 0, \quad \|\boldsymbol{J}\widehat{\boldsymbol{\theta}}^*\|_{\ell_2} - \|\boldsymbol{J}\widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^*\|_{\ell_2} \xrightarrow{\mathcal{P}} 0,$$

where $\stackrel{\mathcal{P}}{\longrightarrow}$ denotes convergence in probability.

Therefore, the GEP (6.17) follows from combining Propositions 6.8 and 6.10.

Finally, the result of Theorem 4.2 follows by computing $AR_{nl}(\widehat{\theta}_{nl}^*)$ when we send $\zeta \to 0$ after $d \to \infty$. This characterization will be carried out in Step 4 of the proof which will be described in the next section.

In non-adversarial contexts and for the standard risk (a.k.a. the generalization error) GEP has been observed by several previous works (see, e.g. [37, 63, 1, 61, 39, 31, 28, 30] and also [56, 12, 66] in the context of random kernel matrices). In [61] the authors provide a precise characterization of the standard risk for the random features model (in non-adversarial setting) and observed that it corresponds to that of its noisy linear counterpart model. A similar GEP phenomena was conjectured for maximum-margin linear classifiers in binary classification [63]. Subsequently, GEP has been proved for more general settings by [31] and [39] for a teacher-student framework. In [31] the authors show GEP for learning with one-pass

stochastic gradient descent (SGD). The work [39] considers the empirical risk minimization (with all data), which results in complicated correlations between the estimator and the samples, and proves the GEP for these settings. However, we cannot directly apply the result of [39] since it assumes that the activation function is an odd function, thrice continuously differentiable with bounded first three derivatives, which are violated for the ReLU activation. Also, our adversarial loss function has additional terms that are beyond the setting considered in [39]. Nevertheless, our proof of Theorem 6.9 is based on the machinery developed in [39]. Here we use a central limit theorem for weakly correlated random variables proved by [30] to show GEP in the context of adversarial training.

6.4. Analysis of the Gaussian noisy linear model via convex Gaussian minimax framework In our final step, we provide a sharp characterization of the adversarial risk $\mathop{\mathsf{AR}}_{\mathsf{nl}}(\widehat{\boldsymbol{\theta}}_{\mathsf{nl}}^*)$ using the Convex Gaussian Minimax Theorem (CGMT), which is a powerful and tight extension of Gordon's Gaussian process inequality [33] with the presence of convexity. The underlying idea of the CGMT framework dates back to [80, 81, 82] where the constrained LASSO was analyzed in the high signal-to-noise ratio regimes. The seminal work [88, 87] significantly extended these ideas and developed the CGMT framework to precisely characterize the mean-squared errors of regularized M- estimators in high-dimensional linear models.

At a more technical level, the CGMT provides a principled machinery to characterize the asymptotic behavior of certain mini-max optimization problems that are affine in a Gaussian matrix X, namely problems of the form

(6.25)
$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{u}} \quad \boldsymbol{u}^{\mathsf{T}} \boldsymbol{X} \boldsymbol{\theta} + \phi(\boldsymbol{\theta}, \boldsymbol{u})$$

where $\phi(\theta, u)$ is convex in θ and concave in u. The CGMT decouples the above objective into a much simpler Gaussian process with essentially the same limit, yet much easier to analyze:

(6.26)
$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{u}} \quad \|\boldsymbol{\theta}\|_{\ell_2} \boldsymbol{g}^{\mathsf{T}} \boldsymbol{u} + \|\boldsymbol{u}\|_{\ell_2} \boldsymbol{h}^{\mathsf{T}} \boldsymbol{\theta} + \phi(\boldsymbol{\theta}, \boldsymbol{u}),$$

where g and h are independent Gaussian vectors with i.i.d N(0,1) entries. We refer to [88, Theorem 3] for a precise statement on the relation between the optimization (6.25), often referred to as *Primary Optimization (PO)* and (6.26), called *Auxiliary Optimization (AO)*. The next step is to derive the point-wise limit of the AO objective in the large dimension limit and showing that it concentrates around a deterministic function with a small number of scalar variables (called *scalarization* step). By showing that this convergence is uniform over a neighborhood of solution and using convexity-concavity of the function, we obtain a precise characterization of the adversarial risk in terms of the solutions of the corresponding convex-concave (deterministic) optimization (4.5).

Note that although the CGMT is a general machinery, the derivation and the study of the AO problem is entirely problem-specific and is usually rather challenging, often requiring the development of non-trivial probabilistic analysis. In relation to *Approximate message passing (AMP)*, which is another powerful tool for deriving asymptotically exact characterization of high-dimensional estimators (see e.g. [25]), it is worth noting that both of these techniques provide a deterministic equation (called state evolution in the AMP parlance) which describes the large limit behavior of a random system.

The CGMT has been recently used in several contexts, e.g., to characterize the performance of high-dimensional regularized logistic regression [73], SLOPE estimator in sparse linear regression [38], boosting and min ℓ_1 norm classifier [52], multi-class classification [89], and phase retrieval [20]. More closely to our work, the CGMT has been used to study the effect of adversarial training in the context of linear regression [46] and linear classifiers [45, 85]. On a technical side, the CGMT analysis for our current problem is more involved and intricate

than the analysis carried out in [46] for linear regression due to: (i) features f_i in (6.23) being correlated; (ii) the presence of the matrix J in the loss which introduces more interactions among the model parameters.

Acknowledgments The authors thank Alexander Robey for interesting discussions and feedback on an early draft.

Funding The research of H. Hassani is supported by the NSF CAREER award, AFOSR YIP, the Intel Rising Star award, as well as the AI Institute for Learning-Enabled Optimization at Scale (TILOS). A. Javanmard is partially supported by the Sloan Research Fellowship in mathematics, an Adobe Data Science Faculty Research Award, the NSF CAREER Award DMS-1844481 and NSF Award DMS-2311024.

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 [36].

REFERENCES

- [1] ABBASI, E., SALEHI, F. and HASSIBI, B. (2019). Universality in Learning from Linear Measurements. *Advances in Neural Information Processing Systems* **32** 12372–12382.
- [2] BAI, Z. and SILVERSTEIN, J. W. (2010). Spectral analysis of large dimensional random matrices 20. Springer.
- [3] BARTLETT, P. L., MONTANARI, A. and RAKHLIN, A. (2021). Deep learning: a statistical viewpoint. *arXiv* preprint arXiv:2103.09177.
- [4] BELKIN, M., HSU, D., MA, S. and MANDAL, S. (2019). Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences* 116 15849– 15854.
- [5] BELKIN, M., MA, S. and MANDAL, S. (2018). To Understand Deep Learning We Need to Understand Kernel Learning. In *International Conference on Machine Learning* 541–549.
- [6] BIGGIO, B., CORONA, I., MAIORCA, D., NELSON, B., ŠRNDIĆ, N., LASKOV, P., GIACINTO, G. and ROLI, F. (2013). Evasion attacks against machine learning at test time. In *Joint European conference on machine learning and knowledge discovery in databases* 387–402. Springer.
- [7] BILLINGSLEY, P. (1995). Probability and Measure. Wiley Series in Probability and Statistics. Wiley.
- [8] BOYD, S. and VANDENBERGHE, L. (2009). Convex optimization. Cambridge university press.
- [9] BUBECK, S., LI, Y. and NAGARAJ, D. M. (2021). A law of robustness for two-layers neural networks. In *Conference on Learning Theory* 804–820. PMLR.
- [10] BUBECK, S. and SELLKE, M. (2021). A Universal Law of Robustness via Isoperimetry. Advances in Neural Information Processing Systems 34.
- [11] CARMON, Y., RAGHUNATHAN, A., SCHMIDT, L., LIANG, P. and DUCHI, J. C. (2019). Unlabeled data improves adversarial robustness. *arXiv* preprint arXiv:1905.13736.
- [12] CHENG, X. and SINGER, A. (2013). The spectrum of random inner-product kernel matrices. *Random Matrices: Theory and Applications* **2** 1350010.
- [13] COHEN, J., ROSENFELD, E. and KOLTER, Z. (2019). Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning* 1310–1320. PMLR.
- [14] DANIELY, A. (2017). SGD learns the conjugate kernel class of the network. In Advances in Neural Information Processing Systems 2422–2430.
- [15] DANIELY, A., FROSTIG, R. and SINGER, Y. (2016). Toward deeper understanding of neural networks: the power of initialization and a dual view on expressivity. In *Proceedings of the 30th International Conference on Neural Information Processing Systems* 2261–2269.
- [16] DENG, Z., KAMMOUN, A. and THRAMPOULIDIS, C. (2019). A Model of Double Descent for High-dimensional Binary Linear Classification. *arXiv* preprint arXiv:1911.05822.

- [17] DENG, Z., ZHANG, L., GHORBANI, A. and ZOU, J. (2021). Improving adversarial robustness via unlabeled out-of-domain data. In *International Conference on Artificial Intelligence and Statistics* 2845–2853. PMLR.
- [18] DEPERSIN, J. and LECUÉ, G. (2023). On the robustness to adversarial corruption and to heavy-tailed data of the Stahel–Donoho median of means. *Information and Inference: A Journal of the IMA* 12 814–850.
- [19] DHIFALLAH, O. and Lu, Y. M. (2020). A precise performance analysis of learning with random features. arXiv preprint arXiv:2008.11904.
- [20] DHIFALLAH, O., THRAMPOULIDIS, C. and LU, Y. M. (2018). Phase retrieval via polytope optimization: Geometry, phase transitions, and new algorithms. *arXiv preprint arXiv:1805.09555*.
- [21] DOBRIBAN, E., HASSANI, H., HONG, D. and ROBEY, A. (2020). Provable tradeoffs in adversarially robust classification. *arXiv preprint arXiv:2006.05161*.
- [22] DOHMATOB, E. (2021). Fundamental tradeoffs between memorization and robustness in random features and neural tangent regimes. *arXiv* preprint arXiv:2106.02630.
- [23] DONHAUSER, K., TIFREA, A., AERNI, M., HECKEL, R. and YANG, F. (2021). Interpolation can hurt robust generalization even when there is no noise. *Advances in Neural Information Processing Systems* 34.
- [24] DONOHO, D. L., MALEKI, A. and MONTANARI, A. (2009). Message-passing algorithms for compressed sensing. Proceedings of the National Academy of Sciences 106 18914–18919.
- [25] DONOHO, D. L., MALEKI, A. and MONTANARI, A. (2009). Message-passing algorithms for compressed sensing. Proceedings of the National Academy of Sciences 106 18914–18919.
- [26] EL KAROUI, N. (2010). The spectrum of kernel random matrices. The Annals of Statistics 38 1-50.
- [27] EL KAROUI, N. (2018). On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields* 170 95–175
- [28] GERACE, F., LOUREIRO, B., KRZAKALA, F., MÉZARD, M. and ZDEBOROVÁ, L. (2020). Generalisation error in learning with random features and the hidden manifold model. In *International Conference on Machine Learning* 3452–3462. PMLR.
- [29] GILMER, J., METZ, L., FAGHRI, F., SCHOENHOLZ, S. S., RAGHU, M., WATTENBERG, M. and GOODFELLOW, I. (2018). Adversarial spheres. arXiv preprint arXiv:1801.02774.
- [30] GOLDT, S., LOUREIRO, B., REEVES, G., KRZAKALA, F., MÉZARD, M. and ZDEBOROVÁ, L. (2020). The Gaussian equivalence of generative models for learning with shallow neural networks. arXiv preprint arXiv:2006.14709.
- [31] GOLDT, S., MÉZARD, M., KRZAKALA, F. and ZDEBOROVÁ, L. (2020). Modeling the influence of data structure on learning in neural networks: The hidden manifold model. *Physical Review X* **10** 041044.
- [32] GOODFELLOW, I. J., SHLENS, J. and SZEGEDY, C. (2015). 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.
- [33] GORDON, Y. (1988). On Milman's inequality and random subspaces which escape through a mesh in \mathbb{R}^n . In Geometric aspects of functional analysis 84–106. Springer.
- [34] GOWAL, S., QIN, C., UESATO, J., MANN, T. and KOHLI, P. (2020). Uncovering the Limits of Adversarial Training against Norm-Bounded Adversarial Examples. *arXiv preprint arXiv:2010.03593*.
- [35] GUNASEKAR, S., LEE, J. D., SOUDRY, D. and SREBRO, N. (2018). Implicit Bias of Gradient Descent on Linear Convolutional Networks. In Advances in Neural Information Processing Systems (S. BENGIO, H. WALLACH, H. LAROCHELLE, K. GRAUMAN, N. CESA-BIANCHI and R. GARNETT, eds.) 31 9461–9471. Curran Associates, Inc.
- [36] HASSANI, H. and JAVANMARD, A. (2022). Supplementary material to "The curse of overparametrization in adversarial training: Precise analysis of robust generalization for random features regression".
- [37] HASTIE, T., MONTANARI, A., ROSSET, S. and TIBSHIRANI, R. J. (2022). Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics* **50** 949–986.
- [38] Hu, H. and Lu, Y. M. (2019). Asymptotics and optimal designs of SLOPE for sparse linear regression. In 2019 IEEE International Symposium on Information Theory (ISIT) 375–379. IEEE.
- [39] Hu, H. and Lu, Y. M. (2020). Universality laws for high-dimensional learning with random features. arXiv preprint arXiv:2009.07669.
- [40] HUANG, S.-T. and LEDERER, J. (2023). DeepMoM: Robust Deep Learning With Median-of-Means. *Journal of Computational and Graphical Statistics* **32** 181–195.
- [41] JACOT, A., GABRIEL, F. and HONGLER, C. (2018). Neural Tangent Kernel: Convergence and Generalization in Neural Networks. In *Advances in Neural Information Processing Systems* (S. BENGIO, H. WALLACH, H. LAROCHELLE, K. GRAUMAN, N. CESA-BIANCHI and R. GARNETT, eds.) **31** 8571–8580. Curran Associates, Inc.
- [42] JACOT, A., GABRIEL, F. and HONGLER, C. (2018). Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems* 8571–8580.

- [43] JALAL, A., ILYAS, A., DASKALAKIS, C. and DIMAKIS, A. G. (2017). The robust manifold defense: Adversarial training using generative models. *arXiv preprint arXiv:1712.09196*.
- [44] JAVANMARD, A., MONDELLI, M. and MONTANARI, A. (2020). Analysis of a two-layer neural network via displacement convexity. *The Annals of Statistics* **48** 3619–3642.
- [45] JAVANMARD, A. and SOLTANOLKOTABI, M. (2022). Precise statistical analysis of classification accuracies for adversarial training. *The Annals of Statistics* 50 2127–2156.
- [46] JAVANMARD, A., SOLTANOLKOTABI, M. and HASSANI, H. (2020). Precise tradeoffs in adversarial training for linear regression. In *Conference on Learning Theory* 2034–2078. PMLR.
- [47] KURAKIN, A., GOODFELLOW, I. and BENGIO, S. (2016). Adversarial machine learning at scale. *arXiv* preprint arXiv:1611.01236.
- [48] LAI, L. and BAYRAKTAR, E. (2020). On the adversarial robustness of robust estimators. IEEE Transactions on Information Theory 66 5097–5109.
- [49] LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. Annals of Statistics 1302–1338.
- [50] LEDOUX, M. (2001). The concentration of measure phenomenon. *volume 89 of Mathematical Surveys and Monographs. American Mathematical Society, Providence, RI.*
- [51] LI, Y. and LIANG, Y. (2018). Learning overparameterized neural networks via stochastic gradient descent on structured data. NeurIPS.
- [52] LIANG, T. and SUR, P. (2020). A Precise High-Dimensional Asymptotic Theory for Boosting and Min-L1-Norm Interpolated Classifiers. arXiv preprint arXiv:2002.01586.
- [53] LIESE, F. and MIESCKE, K.-J. (2008). Statistical decision theory. In *Statistical Decision Theory: Estimation, Testing, and Selection* 1–52. Springer.
- [54] LIESE, F. and MIESCKE, K.-J. (2008). Statistical Decision Theory: Estimation, Testing, and Selection. In *Springer Science & Business Media*.
- [55] LINDEBERG, J. W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitsrechnung. Mathematische Zeitschrift 15 211–225.
- [56] LOUART, C., LIAO, Z. and COUILLET, R. (2018). A random matrix approach to neural networks. The Annals of Applied Probability 28 1190–1248.
- [57] MADRY, A., MAKELOV, A., SCHMIDT, L., TSIPRAS, D. and VLADU, A. (2018). 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.
- [58] MADRY, A., MAKELOV, A., SCHMIDT, L., TSIPRAS, D. and VLADU, A. (2018). Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*.
- [59] MAHLOUJIFAR, S., DIOCHNOS, D. I. and MAHMOODY, M. (2019). The curse of concentration in robust learning: Evasion and poisoning attacks from concentration of measure. In *Proceedings of the AAAI* Conference on Artificial Intelligence 33 4536–4543.
- [60] MAHLOUJIFAR, S. and MAHMOODY, M. (2019). Can Adversarially Robust Learning LeverageComputational Hardness? In *Algorithmic Learning Theory* 581–609. PMLR.
- [61] MEI, S. and MONTANARI, A. (2021). The generalization error of random features regression: Precise asymptotics and double descent curve. Communications on Pure and Applied Mathematics, doi.org/10.1002/cpa.22008.
- [62] MIN, Y., CHEN, L. and KARBASI, A. (2021). The curious case of adversarially robust models: More data can help, double descend, or hurt generalization. In *Uncertainty in Artificial Intelligence* 129–139. PMLR.
- [63] MONTANARI, A., RUAN, F., SOHN, Y. and YAN, J. (2019). The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparametrized regime. arXiv preprint arXiv:1911.01544.
- [64] MONTANARI, A., ZHONG, Y. and ZHOU, K. (2021). Tractability from overparametrization: The example of the negative perceptron. *arXiv* preprint arXiv:2110.15824.
- [65] NAJAFI, A., MAEDA, S.-I., KOYAMA, M. and MIYATO, T. (2019). Robustness to adversarial perturbations in learning from incomplete data. *arXiv preprint arXiv:1905.13021*.
- [66] PENNINGTON, J. and WORAH, P. (2019). Nonlinear random matrix theory for deep learning. *Journal of Statistical Mechanics: Theory and Experiment* **2019** 124005.
- [67] RAGHUNATHAN, A., XIE, S. M., YANG, F., DUCHI, J. C. and LIANG, P. (2019). Adversarial Training Can Hurt Generalization. arXiv preprint arXiv:1906.06032.
- [68] RAHIMI, A. and RECHT, B. (2007). Random features for large-scale kernel machines. *Advances in neural information processing systems* **20** 1177–1184.
- [69] RAHIMI, A. and RECHT, B. (2008). Uniform approximation of functions with random bases. In 2008 46th Annual Allerton Conference on Communication, Control, and Computing 555–561. IEEE.
- [70] REBUFFI, S.-A., GOWAL, S., CALIAN, D. A., STIMBERG, F., WILES, O. and MANN, T. A. (2021). Data Augmentation Can Improve Robustness. *Advances in Neural Information Processing Systems* 34.

- [71] RICHARDSON, T. and URBANKE, R. (2008). Modern coding theory. Cambridge university press.
- [72] RUDELSON, M., VERSHYNIN, R. et al. (2013). Hanson-Wright inequality and sub-gaussian concentration. Electronic Communications in Probability 18.
- [73] SALEHI, F., ABBASI, E. and HASSIBI, B. (2019). The Impact of Regularization on High-dimensional Logistic Regression. In *Advances in Neural Information Processing Systems* (H. WALLACH, H. LAROCHELLE, A. BEYGELZIMER, F. D' ALCHÉ-BUC, E. FOX and R. GARNETT, eds.) 32. Curran Associates, Inc.
- [74] SEHWAG, V., MAHLOUJIFAR, S., HANDINA, T., DAI, S., XIANG, C., CHIANG, M. and MITTAL, P. (2021). Improving Adversarial Robustness Using Proxy Distributions. *arXiv* preprint *arXiv*:2104.09425.
- [75] SHAFAHI, A., HUANG, W. R., STUDER, C., FEIZI, S. and GOLDSTEIN, T. (2019). Are adversarial examples inevitable? In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019.
- [76] SION, M. et al. (1958). On general minimax theorems. Pacific Journal of mathematics 8 171–176.
- [77] SOLTANOLKOTABI, M., JAVANMARD, A. and LEE, J. D. (2018). Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory* 65 742–769.
- [78] SONG, M., MONTANARI, A. and NGUYEN, P. (2018). A mean field view of the landscape of two-layers neural networks. In *Proceedings of the National Academy of Sciences* **115** E7665–E7671.
- [79] SOUDRY, D., HOFFER, E., NACSON, M. S., GUNASEKAR, S. and SREBRO, N. (2018). The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research* 19 2822–2878.
- [80] STOJNIC, M. (2013). A framework to characterize performance of LASSO algorithms. arXiv preprint arXiv:1303.7291.
- [81] STOJNIC, M. (2013). Meshes that trap random subspaces. arXiv preprint arXiv:1304.0003.
- [82] STOJNIC, M. (2013). Upper-bounding ℓ_1 -optimization weak thresholds. arXiv preprint arXiv:1303.7289.
- [83] SU, D., ZHANG, H., CHEN, H., YI, J., CHEN, P.-Y. and GAO, Y. (2018). Is Robustness the Cost of Accuracy?—A Comprehensive Study on the Robustness of 18 Deep Image Classification Models. In Proceedings of the European Conference on Computer Vision (ECCV) 631–648.
- [84] SZEGEDY, C., ZAREMBA, W., SUTSKEVER, I., BRUNA, J., ERHAN, D., GOODFELLOW, I. J. and FER-GUS, R. (2014). Intriguing properties of neural networks. ICLR, abs/1312.6199, 2014.
- [85] TAHERI, H., PEDARSANI, R. and THRAMPOULIDIS, C. (2020). Asymptotic behavior of adversarial training in binary classification. arXiv preprint arXiv:2010.13275.
- [86] THRAMPOULIDIS, C., ABBASI, E. and HASSIBI, B. (2015). Precise high-dimensional error analysis of regularized m-estimators. In 2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton) 410–417. IEEE.
- [87] THRAMPOULIDIS, C., ABBASI, E. and HASSIBI, B. (2018). Precise Error Analysis of Regularized *M*-Estimators in High Dimensions. *IEEE Transactions on Information Theory* **64** 5592–5628.
- [88] THRAMPOULIDIS, C., OYMAK, S. and HASSIBI, B. (2015). Regularized linear regression: A precise analysis of the estimation error. In *Conference on Learning Theory* 1683–1709.
- [89] THRAMPOULIDIS, C., OYMAK, S. and SOLTANOLKOTABI, M. (2020). Theoretical insights into multiclass classification: A high-dimensional asymptotic view. *arXiv preprint arXiv:2011.07729*.
- [90] TSIPRAS, D., SANTURKAR, S., ENGSTROM, L., TURNER, A. and MADRY, A. (2019). Robustness May Be at Odds with Accuracy. In *International Conference on Learning Representations*.
- [91] VERSHYNIN, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint* arXiv:1011.3027.
- [92] VERSHYNIN, R. (2018). *High-dimensional probability: An introduction with applications in data science* **47**. Cambridge University Press.
- [93] WONG, E. and KOLTER, J. Z. (2018). 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 5283–5292.
- [94] Wu, B., Chen, J., Cai, D., He, X. and Gu, Q. (2021). Do Wider Neural Networks Really Help Adversarial Robustness? *Advances in Neural Information Processing Systems* **34**.
- [95] ZHAI, R., CAI, T., HE, D., DAN, C., HE, K., HOPCROFT, J. and WANG, L. (2019). Adversarially robust generalization just requires more unlabeled data. *arXiv preprint arXiv:1906.00555*.
- [96] ZHANG, H., Wu, Y. and HUANG, H. (2022). How Many Data Are Needed for Robust Learning? arXiv preprint arXiv:2202.11592.
- [97] ZHANG, H., YU, Y., JIAO, J., XING, E. P., GHAOUI, L. E. and JORDAN, M. I. (2019). 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 7472– 7482.

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

BY HAMED HASSANI^{1,a}, ADEL JAVANMARD^{2,b}

The supplementary materials contain the proofs of theorems and technical lemmas. It is structured around the main four steps outlined in Section 6.

For the sake of completeness, we reintroduce the notation used throughout the proofs. **Notations.** Throughout the paper, we use $O_d(\cdot)$, $o_d(\cdot)$ to denote the standard big-O and little-o notation, where we stress the asymptotic variable d. Likewise, we denote by $O_{d,\mathbb{P}}$ and $o_{d,\mathbb{P}}$ to indicate asymptotic behavior in probability. Specifically, $f(d) = O_{d,\mathbb{P}}(g(d))$ if for any $\varepsilon > 0$, there exists $C_{\varepsilon} > 0$ and large enough d_{ε} such that $\mathbb{P}(|f(d)/g(d)| > C_{\varepsilon}) \le \varepsilon$, for all $d \ge d_{\varepsilon}$. Similarly, $f(d) = o_{d,\mathbb{P}}(g(d))$ if f(d)/g(d) converges to zero in probability. We write $f(d) \approx g(d)$ as $d \to \infty$, when $f(d) - g(d) \to 0$, in probability. Note that we consider the asymptotic regime where n, d, N grow at the same scale, $(\lim N/d \to \psi_1)$ and $\lim n/d \to \psi_2$ for some positive constants ψ_1 and ψ_2), the expression $d \to \infty$ implies that $n, N \to \infty$, as well.

For a matrix A, we denote by ||A|| its operator norm, $||A||_F = (\sum_{ij} A_{ij}^2)^{1/2}$ the Frobenius norm of A. For two matrices A and B of same size, we let $A \odot B$ be the element-wise product of A and B. In addition, [A; B] concatenates the two matrices row-wise and [A, B] denotes the column-wise concatenation. For an integer n, we use the shorthand $[n] = \{1, \ldots, n\}$.

CONTENTS

A	Inter	changing the limits of $d \to \infty$ and $\zeta \to 0$	
	A. 1	Proof of Lemma A.3	
В	Proo	fs of step 1: Asymptotically-exact closed form of adversarial examples	
	B.1	Proof of Lemma B.1	
	B.2	Proof of Proposition 6.2	
	B.3	Proof of Proposition 6.3	
C		fs of step 2: Concentration of the adversarial effects	
	C.1	Proof of Proposition 6.4	
	C.2	Proof of Lemma C.1	
	C.3	Proof of Lemma C.2	
	C.4	Proof of Lemma 6.6	
	C.5	Proof of Lemma 6.7	
D	Proofs of step 3: The Gaussian equivalence property		
	D.1	Proof of Proposition 6.8	
	D.2	Proof of Theorem 6.9	
	D.3	Proof of Proposition 6.10	
E	Proofs of Step 4: Analysis of the Gaussian noisy linear model via convex Gaussian		
	mini	max framework	
	E.1	Scalarization of the AO problem	
	E.2	Convergence analysis of the AO problem	

1

¹Department of Electrical and Systems Engineering, University of Pennsylvania, ^ahassani@seas.upenn.edu

²Data Sciences and Operations Department, University of Southern California, ^bajavanma@usc.edu

	E.3	Proof of Theorem 4.2(b)	55
	E.4	Proof of Proposition 5.1	57
	E.5	Proofs of the Auxiliary Lemmas	58
7	Som	e useful lemmas	60

APPENDIX A: INTERCHANGING THE LIMITS OF $d \to \infty$ AND $\zeta \to 0$

Consider the loss function (5.1) given by

$$\mathcal{L}(\boldsymbol{\theta}, \zeta, d) = \max_{\|\boldsymbol{\delta}_i\|_{\ell_2} \le \varepsilon} \frac{1}{2n} \sum_{i=1}^n (y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i)))^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \Omega \boldsymbol{\theta},$$

where with a slight abuse of notation, we made the dependence on ζ and d explicit.

In the next lemma we show that the order of the two limits $d \to \infty$ and $\zeta \to 0$ can be interchanged.

Lemma A.1 *Under the assumptions of Theorem 4.2, we have*

$$\lim_{d\to\infty} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, 0, d) = \lim_{\zeta\to 0} \lim_{d\to\infty} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d).$$

Proof (Proof of Lemma A.1) First note that

(A.1)
$$\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, 0, d) = \min_{\boldsymbol{\theta}, \zeta \geq 0} \mathcal{L}(\boldsymbol{\theta}, \zeta, d) = \min_{\zeta \geq 0} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d) = \lim_{\zeta \to 0} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d).$$

The last step holds since $\mathcal{L}(\theta, \zeta, d)$ is increasing in ζ for all θ :

$$\mathcal{L}(\boldsymbol{\theta}, \zeta_1, d) \leq \mathcal{L}(\boldsymbol{\theta}, \zeta_2, d), \quad \text{if } \zeta_1 \leq \zeta_2.$$

Minimizing both sides over θ , we get that $\min_{\theta} \mathcal{L}(\theta, \zeta, d)$ is increasing in ζ . We next show that

(A.2)
$$\lim_{d\to\infty} \lim_{\zeta\to 0} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta},\zeta,d) = \lim_{\zeta\to 0} \lim_{d\to\infty} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta},\zeta,d),$$

where the limits are in probability. Without loss of generality we restrict the domain of ζ to $[0, \zeta_*]$, for an arbitrary but fixed ζ_* . The reason is that in our proofs provided in the paper we allow ζ to be an arbitrarily small fixed value (i.e. we need ζ to be arbitrarily small, but fixed). We next use the Moore-Osgood theorem on exchanging limits, by which we need to verify that

(A.3)
$$\lim_{d\to\infty} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d) = f(\zeta), \quad \text{uniformly on } \zeta \in (0, \zeta_*],$$

(A.4)
$$\lim_{\zeta \to 0} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d) = A_d, \quad \text{pointwise over } d \in \mathbb{N}.$$

The second identity follows from (A.1). To prove the first identity, note that $\mathcal{L}(\theta, \zeta, d)$ is convex in (θ, ζ) . Now, since partial minimization preserves convexity [8, Section 3.2.5], $\min_{\theta} \mathcal{L}(\theta, \zeta, d)$ is convex in ζ . The point-wise limit of (A.3) is already established in the paper, and we obtain uniform convergence using the convexity lemma [53, Lemma 7.75]. In words, the lemma states that pointwise convergence of convex functions implies uniform convergence in compact subsets.

Combining (A.1) and (A.2) we obtain

(A.5)
$$\lim_{d\to\infty} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, 0, d) = \lim_{\zeta\to 0} \lim_{d\to\infty} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d).$$

Recall that our main goal in the paper is to characterize the in-probability limit of $AR(\widehat{\theta}^{\varepsilon})$, with

(A.6)
$$\widehat{\boldsymbol{\theta}}^{\varepsilon} = \arg\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, 0, d).$$

Define

(A.7)
$$\widehat{\boldsymbol{\theta}_{\zeta}^{\varepsilon}} = \arg\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d).$$

Our next lemma relates $AR(\widehat{\boldsymbol{\theta}}^{\varepsilon})$ to $AR(\widehat{\boldsymbol{\theta}}^{\varepsilon})$.

Proposition A.2 Let $\widehat{\theta}^{\varepsilon}$ and $\widehat{\theta}_{\zeta}^{\varepsilon}$ be respectively given by (A.6) and (A.7). Under assumptions of Theorem 4.2, we have

$$\lim_{\zeta \to 0} \lim_{d \to \infty} \mathsf{AR}(\widehat{\boldsymbol{\theta}}_{\zeta}^{\varepsilon}) = \lim_{d \to \infty} \mathsf{AR}(\widehat{\boldsymbol{\theta}}^{\varepsilon}).$$

Proof (Proof of Proposition A.2) To proof the claim, we use a standard trick to translate the question on the optimal solution of the minimization problem (i.e. $\widehat{\theta}^{\varepsilon}$, $\widehat{\theta}^{\varepsilon}_{\zeta}$) to one regarding the optimal costs.

Let $B = \lim_{\zeta \to 0} \lim_{d \to \infty} \mathsf{AR}(\widehat{\boldsymbol{\theta}}_{\zeta}^{\varepsilon})$ (which exists and is calculated in Theorem 4.2). We need to show that $\widehat{\boldsymbol{\theta}}^{\varepsilon}$ belongs to the following set $S_{\delta} \coloneqq \{\boldsymbol{\theta} : |\mathsf{AR}(\boldsymbol{\theta}) - B| \le \delta\}$, with probability converging to one (as $d \to \infty$) for all $\delta > 0$. Let S_{δ}^{ε} denotes the complement set. If we show that

(A.8)
$$\min_{\boldsymbol{\theta} \in S_{\delta}^{c}} \mathcal{L}(\boldsymbol{\theta}, 0, d) > \mathcal{L}(\widehat{\boldsymbol{\theta}}^{\varepsilon}, 0, d),$$

then $\widehat{\theta}^{\varepsilon}$ must lie in S_{δ} . We formalize it in the next lemma.

Lemma A.3 Suppose that there exist constants ℓ , $\tilde{\ell}$ and $\eta > 0$ such that

- $\tilde{\ell} \ge \ell + 2\eta$,
- $\mathcal{L}(\widehat{\boldsymbol{\theta}}^{\varepsilon}, 0, d) < \ell + \eta$ with probability at least 1 p,
- $\min_{\theta \in S_{\delta}^{c}} \mathcal{L}(\theta, 0, d) > \ell \eta$ with probability at least 1 p.

Then,
$$\mathbb{P}(\widehat{\boldsymbol{\theta}}^{\varepsilon} \in S_{\delta}) \geq 1 - 2p$$
.

We then have the following corollary.

Corollary A.4 Suppose that there exist constants $\ell < \tilde{\ell}$ such that $\mathcal{L}(\widehat{\boldsymbol{\theta}}^{\varepsilon}, 0, d) \stackrel{p}{\to} \ell$ and $\min_{\boldsymbol{\theta} \in S^{\varepsilon}} \mathcal{L}(\boldsymbol{\theta}, 0, d) \stackrel{p}{\to} \tilde{\ell}$. Then, $\lim_{d \to \infty} \mathbb{P}(\widehat{\boldsymbol{\theta}}^{\varepsilon} \in S_{\delta}) = 1$, for every $\delta > 0$.

In light of the above corollary, we compare the converging limits: Let $\ell := \lim_{d \to \infty} \mathcal{L}(\widehat{\boldsymbol{\theta}}^{\varepsilon}, 0, d)$ and $\tilde{\ell} := \lim_{d \to \infty} \min_{\boldsymbol{\theta} \in S^{\varepsilon}_{\delta}} \mathcal{L}(\boldsymbol{\theta}, 0, d)$. We need to show $\ell < \tilde{\ell}$. A similar trick has been used in [86, Theorem 6.1 (iii)]. We next use (A.5), by which

$$\ell = \lim_{\zeta \to 0} \lim_{d \to \infty} \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d).$$

By a similar argument,

$$\tilde{\ell} = \lim_{\zeta \to 0} \lim_{d \to \infty} \min_{\boldsymbol{\theta} \in S_{\delta}^{c}} \mathcal{L}(\boldsymbol{\theta}, \zeta, d),$$

where the difference is the domain over which we optimize. Note that in Theorem 4.2 we calculate ℓ as the optimal value of a deterministic convex-concave optimization problem and show that it has a unique solution. Likewise, one can obtain a similar optimization for $\tilde{\ell}$ with the difference that its variables come from a restricted domain, which excludes the optimal solution of the former. By uniqueness of the solution, we conclude that $\ell < \tilde{\ell}$, which completes the proof of proposition.

A.1. Proof of Lemma A.3 Define the event

$$\mathcal{E} \coloneqq \{ \min_{\boldsymbol{\theta} \in S_{\delta}^{c}} \mathcal{L}(\boldsymbol{\theta}, 0, d) > \tilde{\ell} - \eta, \mathcal{L}(\widehat{\boldsymbol{\theta}}^{\varepsilon}, 0, d) < \ell + \eta \}.$$

On this event, using the first condition, we see that (A.8) holds and so $\widehat{\theta}^{\varepsilon} \in S_{\delta}$. So we need to show that $\mathbb{P}(\mathcal{E}) \ge 1 - 2p$, which follows easily from union bounding and using the second and third conditions.

Remark A.2 In [61] the authors derive a precise characterization of the generalization of random features regression in a non-adversarial setting. This work makes a conjecture (see Remark 1 therein) that the generalization error of ridgeless estimator is the same as the min-norm least square estimator. This conjecture amounts to showing that the limits $\lambda \to 0$ (λ the ridge penalty parameter) and $d \to \infty$ can be exchanged. We believe that this conjecture can be proved by following a similar argument of the proof of Lemma A.1.

APPENDIX B: PROOFS OF STEP 1: ASYMPTOTICALLY-EXACT CLOSED FORM OF ADVERSARIAL EXAMPLES

Recall that $x \sim N(0, I_d)$. In the following we will show that conditioned on the event \mathcal{E}_W , defined in (6.2), with probability at least $1 - c/(\log(d))^2 - N^2 d^{-C}$ over the choice of x, we have

(B.1)
$$\sup_{\boldsymbol{\theta} \in \mathsf{C}_{\boldsymbol{\theta}}, \|\boldsymbol{\delta}\|_{\ell_2} \le \varepsilon} \left| \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}(\boldsymbol{x} + \boldsymbol{\delta})) - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}\boldsymbol{x}) - \boldsymbol{\theta}^\mathsf{T} \mathrm{diag}(\mathbb{I}(\boldsymbol{W}\boldsymbol{x} > 0)) \boldsymbol{W} \boldsymbol{\delta} \right| = C \frac{\log(d)}{d^{\frac{1}{6}}}.$$

As a result, we can write

(B.2)

$$\max_{\|\boldsymbol{\delta}\|_{\ell_2} \leq \varepsilon} \left| y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}(\boldsymbol{x} + \boldsymbol{\delta})) \right| = \max_{\|\boldsymbol{\delta}\|_{\ell_2} \leq \varepsilon} \left| y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}) - \langle \boldsymbol{W}^\mathsf{T} \mathrm{diag}(\mathbb{I}(\boldsymbol{W} \boldsymbol{x} > 0)) \boldsymbol{\theta}, \boldsymbol{\delta} \rangle \right| + C \frac{\log(d)}{d^{\frac{1}{6}}},$$

for an absolute constant C > 0, uniformly over $\theta \in C_{\theta}$. The maximization problem in the right-hand side of the above relation has a closed-form solution:

(B.3)
$$\boldsymbol{\delta} = \varepsilon \operatorname{sign}(y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x})) \frac{\boldsymbol{W}^{\mathsf{T}} \operatorname{diag}\mathbb{I}(\boldsymbol{W} \boldsymbol{x} > 0) \boldsymbol{\theta}}{\|\boldsymbol{W}^{\mathsf{T}} \operatorname{diag}\mathbb{I}(\boldsymbol{W} \boldsymbol{x} > 0) \boldsymbol{\theta}\|_{\ell_{\alpha}}},$$

which gives us the desired result (6.4). It thus remains to prove (B.1).

Denote the rows of matrix W by $\{w_1, \dots, w_N\}$ with $w_\ell \in \mathbb{R}^d$. Given x, we define the three sets

$$A(\boldsymbol{x}) = \left\{ \ell : \langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle > d^{-\frac{1}{3}} \right\},$$

$$B(\boldsymbol{x}) = \left\{ \ell : \langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle < -d^{-\frac{1}{3}} \right\},$$

 $C(\boldsymbol{x}) = \left\{ \ell : \langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle \in \left[-d^{-\frac{1}{3}}, d^{-\frac{1}{3}} \right] \right\}.$

We first need to bound the cardinality of the set C(x).

Lemma B.1 With probability $1 - c/(\log(d))^2 - N^2 d^{-C}$ we have

$$|C(\boldsymbol{x})| \le Cd^{\frac{2}{3}}\log(d).$$

The proof of this lemma is given in Section B.1.

Now, for a vector δ we define:

$$\Delta A(\boldsymbol{x}, \boldsymbol{\delta}) = \bigg\{ \ell : \ell \in A(\boldsymbol{x}) \text{ and } \langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle < 0 \bigg\},$$

$$\Delta B(\boldsymbol{x}, \boldsymbol{\delta}) = \left\{ \ell : \ell \in B(\boldsymbol{x}) \text{ and } \langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle > 0 \right\}.$$

In other words, the set $\Delta A(x, \delta)$ (respectively $\Delta B(x, \delta)$) contains all the indices ℓ in A(x) (respectively B(x)) in which the sign of $\langle w_\ell, x \rangle$ and $\langle w_\ell, x + \delta \rangle$ are different. We now prove that when $\|\delta\|_{\ell_2} \leq \varepsilon$ we have

(B.4)
$$|\Delta A(\boldsymbol{x}, \boldsymbol{\delta})|, |\Delta B(\boldsymbol{x}, \boldsymbol{\delta})| \le Cd^{\frac{2}{3}},$$

for an absolute constant C>0. By definition, on the event $\mathcal{E}_{\boldsymbol{W}}$, \boldsymbol{W} has bounded operator norm, say at most C for some constant C>0. In addition, $\|\boldsymbol{\delta}\|_{\ell_2} \leq \varepsilon$. Therefore, $\|\boldsymbol{W}\boldsymbol{\delta}\|_{\ell_2} \leq \varepsilon C$. As a result, the number of entries of the vector $\boldsymbol{W}\boldsymbol{\delta}$ whose absolute value is larger than $d^{-\frac{1}{3}}$ is bounded by $\varepsilon^2 C^2 d^{\frac{2}{3}}$. But from the definitions of the sets $A(\boldsymbol{x})$ and $\Delta A(\boldsymbol{x},\boldsymbol{\delta})$ it is immediate that for $\ell \in \Delta A(\boldsymbol{x},\boldsymbol{\delta})$ we have $|\langle \boldsymbol{w}_\ell,\boldsymbol{\delta}\rangle| > d^{-\frac{1}{3}}$. And this results in the fact that $|\Delta A(\boldsymbol{x},\boldsymbol{\delta})| \leq C' d^{\frac{2}{3}}$ with $C' = \varepsilon^2 C^2$. The same argument holds for $|\Delta B(\boldsymbol{x},\boldsymbol{\delta})|$.

Let us now consider the two vectors $\sigma(\boldsymbol{W}\boldsymbol{x})$ and $\sigma(\boldsymbol{W}(\boldsymbol{x}+\boldsymbol{\delta}))$. We would like to find out how these vectors are different on the indices that belong to the set $A(\boldsymbol{x})$ or $B(\boldsymbol{x})$. Let us start with the indices in $B(\boldsymbol{x})$. Note that for any $\ell \in B(\boldsymbol{x})$ we have $\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle < 0$. For this entry, it is easy to see that the two vectors $\sigma(\boldsymbol{W}\boldsymbol{x})$ and $\sigma(\boldsymbol{W}(\boldsymbol{x}+\boldsymbol{\delta}))$ take different values only if we also have $\ell \in \Delta B(\boldsymbol{x}, \boldsymbol{\delta})$. As a result, we can conclude that the two vectors $\sigma(\boldsymbol{W}\boldsymbol{x})$ and $\sigma(\boldsymbol{W}(\boldsymbol{x}+\boldsymbol{\delta}))$ are the same on all the indices belonging to the set $B(\boldsymbol{x})$ except at most $Cd^{\frac{2}{3}}$ indices. In other words, the difference $\sigma(\boldsymbol{W}(\boldsymbol{x}+\boldsymbol{\delta})) - \sigma(\boldsymbol{W}\boldsymbol{x})$ takes zero value on all the indices belonging to the set $B(\boldsymbol{x})$ except at most $Cd^{\frac{2}{3}}$ indices.

For the indices in the set A(x) the situation is different as we are operating in the non-constant part of the ReLU function (note that for any $\ell \in A(x)$ we have $\langle w_\ell, x \rangle > 0$). We first claim the following: The two vectors $\sigma(W(x+\delta))$ and $W(x+\delta)-1/\sqrt{2\pi}$ are the same on all the entries in the set A(x) except the indices in the set $\Delta A(x,\delta)$. The justification is as follows: Consider an index $\ell \in A(x)$ such that $\sigma(\langle w_\ell, x+\delta \rangle) \neq \langle w_\ell, x+\delta \rangle - 1/\sqrt{2\pi}$. Since $\ell \in A(x)$, we have $\sigma(\langle w_\ell, x \rangle) = \langle w_\ell, x \rangle - 1/\sqrt{2\pi}$. Now, since $\sigma(\langle w_\ell, x+\delta \rangle) \neq \langle w_\ell, x+\delta \rangle - 1/\sqrt{2\pi}$ only if $\langle w_\ell, x+\delta \rangle < 0$, we obtain that $\sigma(\langle w_\ell, x+\delta \rangle) \neq \langle w_\ell, x+\delta \rangle - 1/\sqrt{2\pi}$ only if $\ell \in \Delta A(x,\delta)$.

In summary, we have shown that (i) On indices belonging to the set $A(x) \setminus \Delta A(x, \delta)$: the two vectors $\sigma(W(x+\delta))$ and $W(x+\delta) - 1/\sqrt{2\pi}$ are the same, and (ii) on the indices belonging to the set $B(x) \setminus \Delta B(x, \delta)$ the vector $\sigma(W(x+\delta)) - \sigma(Wx)$ takes value 0. Also, (iii) both sets $\Delta A(x, \delta)$ and $\Delta B(x, \delta)$ have cardinality at most $Cd^{\frac{2}{3}}$. We can thus write:

$$(B.5) = \sum_{\ell \in A(\boldsymbol{x})} \theta_{\ell}(\sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle) - \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle)) + \sum_{\ell \in B(\boldsymbol{x})} \theta_{\ell}(\sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle) - \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle)) + \sum_{\ell \in B(\boldsymbol{x})} \theta_{\ell}(\sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle) - \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle)).$$

We first bound the second and third terms. Using the fact that $|\sigma(a+b) - \sigma(b)| \le |a|$ we obtain for the third term that:

$$\Big| \sum_{\ell \in C(\boldsymbol{x})} \theta_{\ell} (\sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle) - \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle)) \Big| \leq \sum_{\ell \in C(\boldsymbol{x})} |\theta_{\ell}| |\langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle|$$

$$\leq \|\boldsymbol{\theta}\|_{\infty} \sqrt{|C(\boldsymbol{x})|} \|\boldsymbol{W}\boldsymbol{\delta}\|_{\ell_{2}}$$

$$\leq \frac{C}{d^{\frac{1}{2}}} (Cd^{\frac{2}{3}} \log(d))^{\frac{1}{2}} C\varepsilon$$

$$\leq C' d^{-\frac{1}{6}} \log(d),$$
(B.6)

where we have use the result of Lemma B.1 to bound the size of the set C(x) (and hence the above holds with the probability given in that lemma). For the second term we have (B.7)

$$\Big| \sum_{\ell \in B(\boldsymbol{x})} \theta_{\ell}(\sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle) - \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle)) \Big| \leq \sum_{\ell \in \Delta B(\boldsymbol{x}, \boldsymbol{\delta})} |\theta_{\ell}| |\langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle| \leq ||\boldsymbol{\theta}||_{\infty} \sqrt{|\Delta B(\boldsymbol{x}, \boldsymbol{\delta})|} ||\boldsymbol{W}\boldsymbol{\delta}||_{\ell_{2}} \leq C' d^{-\frac{1}{6}}.$$

Finally, in a similar manner as above we can bound

(B.8)
$$\left| \sum_{\ell \in \Delta A(\boldsymbol{x},\boldsymbol{\delta})} \theta_{\ell}(\sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle) - \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle) \right| \leq \|\boldsymbol{\theta}\|_{\infty} \sqrt{|\Delta A(\boldsymbol{x},\boldsymbol{\delta})|} \|\boldsymbol{W}\boldsymbol{\delta}\|_{\ell_{2}} \leq C' d^{-\frac{1}{6}}.$$

By plugging (B.6), (B.7), and (B.8) into (B.5) we have shown that

$$\boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x} + \boldsymbol{\delta})) - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}) = \sum_{\ell \in A(\boldsymbol{x}) \setminus \Delta A(\boldsymbol{x}, \boldsymbol{\delta})} \theta_{\ell} (\sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle) - \sigma(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle)) + C' d^{-\frac{1}{6}} \log(d)$$

$$= \sum_{\ell \in A(\boldsymbol{x}) \setminus \Delta A(\boldsymbol{x}, \boldsymbol{\delta})} \theta_{\ell} (\langle \boldsymbol{w}_{\ell}, \boldsymbol{x} + \boldsymbol{\delta} \rangle - \langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle) + C' d^{-\frac{1}{6}} \log(d)$$

$$= \sum_{\ell \in A(\boldsymbol{x}) \setminus \Delta A(\boldsymbol{x}, \boldsymbol{\delta})} \theta_{\ell} \langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle + C' d^{-\frac{1}{6}} \log(d),$$
(B.9)

where the second equality follows form the definition of the set $\Delta A(x, \delta)$.

As a final step, we define the set $A^+(x) = \{\ell : \langle w_\ell, x \rangle > 0\}$. Note that $A(x) \subseteq A^+(x)$ and $A^+(x) \setminus A(x) \subseteq C(x)$. As a result $|A^+(x) \setminus A(x)| \le Cd^{\frac{2}{3}} \log(d)$. We thus obtain

$$(B.10) \sum_{\ell \in A(\boldsymbol{x}) \backslash \Delta A(\boldsymbol{x}, \boldsymbol{\delta})} \theta_{\ell} \langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle = \sum_{\ell \in A^{+}(\boldsymbol{x})} \theta_{\ell} \langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle - \sum_{\ell \in (A^{+}(\boldsymbol{x}) \backslash A(\boldsymbol{x})) \cup \Delta A(\boldsymbol{x}, \boldsymbol{\delta})} \theta_{\ell} \langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle$$

$$= \sum_{\ell \in A^{+}(\boldsymbol{x})} \theta_{\ell} \langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle + C' d^{-\frac{1}{6}} \log(d),$$

where the last relations follows from the fact that $|(A^+(x)\backslash A(x)) \cup \Delta A(x,\delta)| \le |C(x)| + |\Delta A(x,\delta)| = O(d^{\frac{2}{3}}\log(d))$. By plugging (B.10) into (B.9) we have

$$\boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x} + \boldsymbol{\delta})) - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}) = \sum_{\ell \in A^{+}(\boldsymbol{x})} \theta_{\ell} \langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle + C d^{-\frac{1}{6}} \log(d)$$

$$= \sum_{\ell : \langle \boldsymbol{w}_{\ell}, \boldsymbol{x} \rangle > 0} \theta_{\ell} \langle \boldsymbol{w}_{\ell}, \boldsymbol{\delta} \rangle + C' d^{-\frac{1}{6}} \log(d)$$

$$= \boldsymbol{\theta}^{\mathsf{T}} \operatorname{diag} \mathbb{I}(\boldsymbol{W} \boldsymbol{x} > 0) \boldsymbol{W} \boldsymbol{\delta} + C' d^{-\frac{1}{6}} \log(d),$$

which is the result of (B.1).

B.1. Proof of Lemma B.1 Define the random variables $\mu_{\ell} := \boldsymbol{w}_{\ell}^{\mathsf{T}} \boldsymbol{x}$ for $\ell = 1, \ldots, N$. Note that since \boldsymbol{x} is gaussian and $\|\boldsymbol{w}_{\ell}\|_{\ell_{2}} = 1$, then $\mu_{\ell} \sim \mathsf{N}(0,1)$. Also, note that μ_{ℓ} 's are correlated with each other and each pair (μ_{ℓ}, μ_{k}) is a jointly-gaussian random variable with correlation $\rho_{\ell,k} := \mathbb{E}[\mu_{\ell}, \mu_{k}] = \boldsymbol{w}_{\ell}^{\mathsf{T}} \boldsymbol{w}_{k}$. Define $z_{\ell} = \mathbf{1}\{\mu_{\ell} \in [-d^{-\frac{1}{3}}, d^{-\frac{1}{3}}]\}$. Note that $\mathbb{E}[z_{\ell}] \leq cd^{-\frac{1}{3}}$.

Define the event $\mathcal{E} \coloneqq \{|\boldsymbol{w}_{\ell}^{\mathsf{T}}\boldsymbol{w}_{k}| \le d^{-1/2}\sqrt{C\log(d)}, \ \forall \ell, k \in [N]\}$. Since $\boldsymbol{w}_{\ell} \sim_{i.i.d} \mathrm{Unif}(\mathbb{S}^{\mathrm{d}-1})$, it is easy to see that $\mathbb{P}(\mathcal{E}) \ge 1 - N^{2}d^{-C}$. We also have

$$\mathbb{P}(\mu_{\ell} = 1, \mu_{k} = 1) = \int_{\mu \in [-d^{-\frac{1}{3}}, d^{-\frac{1}{3}}]} f(\mu_{\ell} = \mu) \mathbb{P}(\mu_{k} \in [-d^{-\frac{1}{3}}, d^{-\frac{1}{3}}] \mid \mu_{\ell} = \mu) d\mu,$$

where f denotes pdf of μ_{ℓ} . Now, given $\mu_{\ell} = \mu$, the distribution of μ_{k} is $N(\mu \rho_{\ell,k}, (1 - \rho_{\ell,k}^{2}))$. It is easy to see that on the event \mathcal{E} , for $|\mu| \leq d^{-\frac{1}{3}}$ we have

$$\mathbb{P}(\mu_k \in [-d^{-\frac{1}{3}}, d^{-\frac{1}{3}}] \mid \mu_\ell = \mu) \le cd^{-\frac{1}{3}},$$

and thus

$$\mathbb{E}[z_{\ell}z_{k}] = \mathbb{P}(\mu_{\ell} = 1, \mu_{k} = 1) \le cd^{-\frac{1}{3}} \int_{\mu \in [-d^{-\frac{1}{3}}, d^{-\frac{1}{3}}]} f(\mu_{\ell} = \mu) \mathrm{d}\mu \le cd^{-\frac{2}{3}}.$$

Let us now consider the average $\bar{z} = \frac{1}{N} \sum_{\ell=1}^{N} z_{\ell}$. We have

$$\mathbb{E}[|\bar{z} - \mathbb{E}[\bar{z}]|^2] = \mathbb{E}[\bar{z}^2] - \mathbb{E}[\bar{z}]^2 \le \mathbb{E}[\bar{z}^2] \le \frac{1}{N^2} \sum_{\ell=1}^{N} \mathbb{E}[z_\ell z_k] \le cd^{-\frac{2}{3}},$$

and thus we obtain via the Chebyshev's inequality that

$$\mathbb{P}\left\{|\bar{z} - \mathbb{E}[\bar{z}]| \ge d^{-\frac{1}{3}}\log(d); \mathcal{E}\right\} \le \frac{c}{(\log(d))^2}.$$

Now, by noticing that $|C(x)|/N = \bar{z}$, and $\mathbb{E}[\bar{z}] \le cd^{-\frac{1}{3}}$, along with the assumption that N, d grows proportionally we obtain

$$\mathbb{P}\left\{|C(\boldsymbol{x})| \geq Cd \times d^{-\frac{1}{3}}\log(d); \mathcal{E}\right\} \leq \frac{c}{(\log(d))^2}.$$

Finally, we have

$$\mathbb{P}\left\{\left|C(\boldsymbol{x})\right| \geq Cd^{\frac{2}{3}}\log(d)\right\} \leq \frac{c}{\log(d)^2} + \mathbb{P}(\mathcal{E}^c) \leq \frac{c}{(\log(d))^2} + N^2d^{-C}.$$

B.2. Proof of Proposition 6.2 Recall the loss $\mathcal{L}(\theta)$ given by

$$\mathcal{L}(\boldsymbol{\theta}) \coloneqq \max_{\|\boldsymbol{\delta}_i\|_{\ell_2} \le \varepsilon} \frac{1}{2n} \sum_{i=1}^n (y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i))^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta},$$

and $\widehat{\boldsymbol{\theta}} = \arg\min \mathcal{L}(\boldsymbol{\theta})$. Bounding $\|\widehat{\boldsymbol{\theta}}\|_{\ell_2}$ is straightforward. By optimality of $\widehat{\boldsymbol{\theta}}$ and comparing the loss at $\widehat{\boldsymbol{\theta}}$ and $\boldsymbol{0}$ we get

$$\frac{\zeta}{2}\widehat{\boldsymbol{\theta}}^{\mathsf{T}}\mathbf{\Omega}\widehat{\boldsymbol{\theta}} \leq \mathcal{L}(\mathbf{0}) = \frac{1}{2n}\sum_{i=1}^{n}y_{i}^{2} < C,$$

with probability at least $1 - e^{-cn}$, for absolute constants c, C > 0. Since $\Omega \ge I$, this implies that $\|\theta\|_{\ell_2} \le C_0$ for sufficiently large C_0 .

To bound $\|\widehat{\theta}\|_{\ell_{\infty}}$ we just need to bound any given entry of $\widehat{\theta}$, e.g. its last entry, with high probability. By symmetry, all the entries have the same marginal distribution. Consequently, each entry of $\widehat{\theta}$ can be analyzed in the same way and $\|\widehat{\theta}\|_{\ell_{\infty}}$ can then be controlled by using the union bound.

With a slight abuse of notation, we consider a (N+1) dimensional version of the above optimization over $[\theta; u]$ and denote the last coordinate of the optimal solution by \hat{u} . Let $\lambda := \frac{\sqrt{\log(d)}}{d}$ and

$$\mathbf{\Omega} = \begin{pmatrix} \widetilde{\mathbf{\Omega}} & \lambda \mathbf{1} \\ \lambda \mathbf{1}^\mathsf{T} & \lambda + 1 \end{pmatrix},$$

where $\widetilde{\Omega}$ is of size N and $\mathbf{1} = (1, 1, ..., 1) \in \mathbb{R}^N$. The last coordinate \hat{u} can be expressed as

$$\hat{u} = \arg\min_{u} \min_{\boldsymbol{\theta}} \left[\frac{1}{2n} \sum_{i=1}^{n} \max_{\|\boldsymbol{\delta}_{i}\|_{\ell_{2}} \leq \varepsilon} \left(y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{*}) - u \sigma(\boldsymbol{w}_{N+1}^{\mathsf{T}}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{*})) \right)^{2} + \frac{\zeta}{2} (\boldsymbol{\theta}^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} \boldsymbol{\theta} + 2\lambda (\mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}) u + (\lambda + 1) u^{2}) \right].$$
(B.11)

We next define f(u) as the objective function of u in (B.11). We proceed by deriving a lower bound for f(u).

Let θ_* be the optimal θ if we set u=0 and denote by $\delta_i^{\setminus u}$ the maximizing δ_i , when u=0. Note that $\delta_i^{\setminus u}$ is in general a function of θ . In addition, define

$$\ell([\boldsymbol{\theta}, u]) = \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i^{\mathsf{V}u}) - u\sigma(\boldsymbol{w}_{N+1}^{\mathsf{T}}(\boldsymbol{x}_i + \boldsymbol{\delta}_i^{\mathsf{V}u})) \right)^2,$$

$$Q([\boldsymbol{\theta}, u]) = \frac{\zeta}{2} (\boldsymbol{\theta}^{\mathsf{T}} \widetilde{\Omega} \boldsymbol{\theta} + 2\lambda (\mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}) u + (\lambda + 1) u^2).$$

Since the pointwise maximum of convex functions is convex, the function $\ell(\cdot)$ is convex and hence we have

(B.12)
$$\ell([\boldsymbol{\theta};u]) \ge \ell([\boldsymbol{\theta}_*;0]) + \langle \nabla_{\boldsymbol{\theta}} \ell([\boldsymbol{\theta},u])|_{[\boldsymbol{\theta}_*;0]}, \boldsymbol{\theta} - \boldsymbol{\theta}_* \rangle + \nabla_u \ell([\boldsymbol{\theta},u])|_{[\boldsymbol{\theta}_*;0]} u.$$
For quadratic function $O([\boldsymbol{\theta},u])$ we have

For quadratic function $Q([\theta; u])$ we have

$$Q([\boldsymbol{\theta}; u]) = \frac{\zeta}{2} \boldsymbol{\theta}_{*}^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} \boldsymbol{\theta}_{*} + \zeta \boldsymbol{\theta}_{*}^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{*}) + \frac{\zeta}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{*})^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{*}) + \frac{\zeta}{2} (2\lambda u \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta} + (\lambda + 1)u^{2})$$

(B.13)
$$= Q([\boldsymbol{\theta}_*; 0]) + \zeta(\boldsymbol{\theta}_*^{\mathsf{T}}\widetilde{\boldsymbol{\Omega}}(\boldsymbol{\theta} - \boldsymbol{\theta}_*) + \lambda(\mathbf{1}^{\mathsf{T}}\boldsymbol{\theta}_*)u) + \frac{\zeta}{2}\{(\boldsymbol{\theta} - \boldsymbol{\theta}_*)^{\mathsf{T}}\widetilde{\boldsymbol{\Omega}}(\boldsymbol{\theta} - \boldsymbol{\theta}_*)$$

(B.14)
$$+ (\lambda + 1)u^2 + 2\lambda u \mathbf{1}^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_*) \}.$$

Combining (B.12) and (B.13) we get

$$\mathcal{L}([\boldsymbol{\theta}; u]) \geq \ell([\boldsymbol{\theta}; u]) + Q([\boldsymbol{\theta}; u])$$

$$\geq \mathcal{L}([\boldsymbol{\theta}_{*}; 0]) + \langle \nabla_{\boldsymbol{\theta}} \ell([\boldsymbol{\theta}, u]) |_{[\boldsymbol{\theta}_{*}; 0]} + \zeta \boldsymbol{\theta}_{*}^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}}, \boldsymbol{\theta} - \boldsymbol{\theta}_{*}) + (\nabla_{u} \ell([\boldsymbol{\theta}, u]) |_{[\boldsymbol{\theta}_{*}; 0]} + \zeta \lambda \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}_{*}) u$$
(B.15)
$$+ \frac{\zeta}{2} \left\{ (\boldsymbol{\theta} - \boldsymbol{\theta}_{*})^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{*}) + (\lambda + 1) u^{2} + 2\lambda u \mathbf{1}^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{*}) \right\}.$$

Here, the first inequality holds since $\mathcal{L}(\cdot)$ involves maximization over δ_i , while in definition of $\ell(\cdot)$ we consider $\delta_i^{\setminus u}$. Though, note that $\mathcal{L}([\theta_*,0]) = \ell([\theta_*;0]) + Q([\theta_*;0])$ because when u=0, $\delta_i^{\setminus u}$ are the maximizing perturbations by definition. We used this observation in the second inequality above.

We argue that the second term in the right-hand side is zero. To see this, first write the partial derivative $\nabla_{\theta} \ell$ as

$$\nabla_{\boldsymbol{\theta}} \ell([\boldsymbol{\theta}, u])$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \left(y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\backslash u}) - u \sigma(\boldsymbol{w}_{N+1}^{\mathsf{T}}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\backslash u})) \right) \times \left(\frac{\partial}{\partial \boldsymbol{\theta}} \left[\boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\backslash u})) + u \frac{\partial}{\partial \boldsymbol{\theta}} \sigma(\boldsymbol{w}_{N+1}^{\mathsf{T}}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\backslash u})) \right) \right).$$

(Note that $\pmb{\delta}_i^{\backslash u}$ is a function of $\pmb{\theta}$.) Therefore,

(B.16)

$$\nabla_{\boldsymbol{\theta}} \ell([\boldsymbol{\theta}, u])|_{[\boldsymbol{\theta}^*; 0]} = -\frac{1}{n} \sum_{i=1}^{n} \left(y_i - \boldsymbol{\theta}_*^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i^{\mathsf{V}u})) \right) \left(\frac{\partial}{\partial \boldsymbol{\theta}} [\boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i^{\mathsf{V}u}))]|_{[\boldsymbol{\theta}_*; 0]} \right).$$

By using the first-order optimality condition for θ_* we have

$$\nabla_{\boldsymbol{\theta}} \ell([\boldsymbol{\theta}, u])|_{[\boldsymbol{\theta}_*; 0]} + \zeta \boldsymbol{\theta}_*^{\mathsf{T}} \widetilde{\Omega}$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \left(y_{i} - \boldsymbol{\theta}_{*}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\mathsf{V}u})) \right) \left(\frac{\partial}{\partial \boldsymbol{\theta}} \left[\boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\mathsf{V}u})) \right] \Big|_{\boldsymbol{\theta}_{*};0]} \right) + \zeta \boldsymbol{\theta}_{*}^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} = 0.$$

Using the above relation in (B.15) we arrive at

$$\mathcal{L}([\boldsymbol{\theta}; u]) \ge \mathcal{L}([\boldsymbol{\theta}_*; 0]) + (\nabla_u \ell([\boldsymbol{\theta}, u])|_{[\boldsymbol{\theta}_*; 0]} + \zeta \lambda \mathbf{1}^\mathsf{T} \boldsymbol{\theta}_*) u$$

$$+ \frac{\zeta}{2} \left\{ (\boldsymbol{\theta} - \boldsymbol{\theta}_*)^\mathsf{T} \widetilde{\boldsymbol{\Omega}} (\boldsymbol{\theta} - \boldsymbol{\theta}_*) + (\lambda + 1) u^2 + 2\lambda u \mathbf{1}^\mathsf{T} (\boldsymbol{\theta} - \boldsymbol{\theta}_*) \right\}.$$

Therefore, by minimizing the both sides over θ we obtain

$$f(u) = \min_{\boldsymbol{\theta}} \mathcal{L}([\boldsymbol{\theta}; u])$$

$$\geq f(0) + (\nabla_{u} \ell([\boldsymbol{\theta}, u])|_{[\boldsymbol{\theta}_{*}; 0]} + \zeta \lambda \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}_{*}) u$$

$$+ \min_{\boldsymbol{\theta}} \frac{\zeta}{2} \left\{ (\boldsymbol{\theta} - \boldsymbol{\theta}_{*})^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{*}) + (\lambda + 1) u^{2} + 2\lambda u \mathbf{1}^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{*}) \right\}$$

$$= f(0) + (\nabla_{u} \ell([\boldsymbol{\theta}, u])|_{[\boldsymbol{\theta}_{*}; 0]} + \lambda \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}_{*}) u + \frac{\zeta}{2} u^{2} (1 + \lambda - \lambda^{2} \mathbf{1}^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}}^{-1} \mathbf{1}).$$
(B.17)

By definition of $\widetilde{\Omega}$, it has 1 as an eigenvector with eigenvalue $1 + \lambda d$. So,

(B.18)
$$1 + \lambda - \lambda^2 \mathbf{1}^\mathsf{T} \widetilde{\mathbf{\Omega}}^{-1} \mathbf{1} \ge 1 + \lambda - \frac{\lambda^2 d}{1 + \lambda d} > 1.$$

By optimality of \hat{u} , we have $f(\hat{u}) \le f(0)$, which together with (B.17) and (B.18) imply that

(B.19)
$$|\hat{u}| \leq \frac{2}{\zeta} |\nabla_{u} \ell([\boldsymbol{\theta}, u])|_{[\boldsymbol{\theta}_{*}; 0]} + \zeta \lambda \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}_{*}|.$$

We next bound the terms on the right-hand side separately. We have

$$\lambda \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}_{\star} \leq \lambda \|\mathbf{1}\|_{\ell_{2}} \|\boldsymbol{\theta}_{\star}\|_{\ell_{2}} \leq \sqrt{\frac{\log(d)}{d}} \|\boldsymbol{\theta}_{\star}\|_{\ell_{2}}$$

By optimality of θ_* (when we set u = 0) and comparing it with $\mathbf{0}$ we get

$$\frac{\zeta}{2} \boldsymbol{\theta}_{\star}^{\mathsf{T}} \widetilde{\boldsymbol{\Omega}} \boldsymbol{\theta}_{\star} \leq \frac{1}{2n} \sum_{i=1}^{n} y_{i}^{2} < C,$$

with probability at least $1 - e^{-cn}$.

Sine $\widetilde{\Omega} \geq I$, this implies that $\lambda \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}_* = O_{\mathbb{P}}(\sqrt{\log(d)/d})$.

To bound the other term, recall that by definition $\frac{\partial}{\partial u} \delta_i^{\setminus u} = 0$ and so $\nabla_u \ell([\theta, u])|_{[\theta_*;0]}$ is given by

(B.20)
$$\nabla_{u}\ell([\boldsymbol{\theta},u])|_{[\boldsymbol{\theta}_{*};0]} = \frac{1}{n} \sum_{i=1}^{n} \left(y_{i} - \boldsymbol{\theta}_{*}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\backslash u})) \sigma(\boldsymbol{w}_{N+1}^{\mathsf{T}}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\backslash u})) \right).$$

To simplify the notation define $m_i := \frac{1}{\sqrt{n}} (y_i - \boldsymbol{\theta}_*^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i^{\mathsf{V}u})))$ and $\boldsymbol{X} = [\boldsymbol{x}_1 | \dots | \boldsymbol{x}_n]^{\mathsf{T}}$. Consider the following event:

$$\mathcal{E}\coloneqq\left\{\left\|\boldsymbol{m}\right\|_{\ell_{2}}\leq C,\;\frac{1}{\sqrt{d}}\left\|\boldsymbol{X}\right\|\leq C\right\}\,,$$

where $m = (m_1, ..., m_n)^T$ and C > 0 is a sufficiently large constant. We show that \mathcal{E} is a high probability event. To see this, first observe that

(B.21)
$$\|\boldsymbol{m}\|_{\ell_{2}}^{2} = \frac{1}{n} \sum_{i=1}^{n} \left(y_{i} - \boldsymbol{\theta}_{*}^{\mathsf{T}} \sigma (\boldsymbol{W}(\boldsymbol{x}_{i} + \boldsymbol{\delta}_{i}^{\mathsf{U}}))^{2} \leq \frac{1}{n} \sum_{i=1}^{n} y_{i}^{2},$$

where the inequality follows from optimality of θ_* and comparing the loss value $\mathcal{L}([\theta_*;0])$ with $\mathcal{L}([0;0])$. Therefore,

$$\mathbb{P}(\|\boldsymbol{m}\|_{\ell_2} > C) \le \mathbb{P}(\frac{1}{\sqrt{n}} \|\boldsymbol{y}\|_{\ell_2} > C) \le e^{-C'n},$$

for absolute constants C, C' (depending on the noise variance τ^2). Also, given that X has i.i.d standard normal entries we have

$$\mathbb{P}(\frac{1}{\sqrt{d}} \|\boldsymbol{X}\| > C) \le 2e^{-cn}.$$

Putting the last two bounds together we obtain $\mathbb{P}(\mathcal{E}^c) \leq 3e^{-cn}$.

Let \mathcal{F} be the σ -algebra generated by an arbitrary $\boldsymbol{X}, \boldsymbol{W}, \boldsymbol{y}$ in \mathcal{E} . Clearly, m_i are measurable with respect to \mathcal{F} . Also, \boldsymbol{w}_{N+1} is drawn independently from \mathcal{F} and hence conditioned on that $\boldsymbol{w}_{N+1}^{\mathsf{T}}\boldsymbol{x}_i \sim \mathsf{N}(0,1)$. Since $\mathbb{E}[\sigma(G)] = 0$ for $G \sim \mathsf{N}(0,1)$, we have $\mathbb{E}[\sigma(\boldsymbol{w}_{N+1}^{\mathsf{T}}\boldsymbol{x}_i)|\mathcal{F}] = 0$. To bound $\nabla_u \ell([\boldsymbol{\theta}, u])|_{[\boldsymbol{\theta}_*;0]}$ we view that as a function of \boldsymbol{w}_{N+1} and condition on \mathcal{F} . We then have

$$\mathbb{E}\left[\nabla_{u}\ell([\boldsymbol{\theta},u])|_{[\boldsymbol{\theta}_{*};0]}|\mathcal{F}\right] = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}m_{i}\,\mathbb{E}\left[\sigma(\boldsymbol{w}_{N+1}^{\mathsf{T}}\boldsymbol{x}_{i})|\mathcal{F}\right] = 0.$$

Also this is a Lipschitz continuous function of w_{N+1} with a Lipschitz factor at most $\frac{C}{\sqrt{d}} \| \mathbf{X} \mathbf{m} \|_{\ell_2} \le \frac{1}{\sqrt{\psi_2}} \frac{1}{\sqrt{d}} \| \mathbf{X} \| \| \mathbf{m} \|_{\ell_2} \le \frac{C^2}{\sqrt{\psi_2}} \coloneqq C_0$. Since w_{N+1} is chosen uniformly at random from the unit sphere, we can apply the concentration bound for Lipschitz function (see e.g. [92, Theorem 5.1.4]), which implies that

(B.22)
$$\mathbb{P}\left(\nabla_{u}\ell([\boldsymbol{\theta},u])|_{[\boldsymbol{\theta}_{*};0]} > t\right) \leq 2e^{-c'dt^{2}}.$$

Choosing $t = C\sqrt{\frac{\log(d)}{d}}$ and using this bound in (B.19) we get

$$|\hat{u}| \le C' \sqrt{\frac{\log(d)}{d}}$$
,

with probability at least $1-4e^{-cn}-2d^{-c'C^2}$. The result follows by choosing C>0 large enough so that $c'C^2>1$ and union bounding over the N coordinates of $\widehat{\theta}$, along with the assumption that N,n,d grow at the same order.

B.3. Proof of Proposition 6.3 We know from the result of Proposition 6.1, or more precisely the equations (B.1)-(B.2) in its proof, that with probability $1 - o_d(1)$ we have

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \left| \max_{\|\boldsymbol{\delta}\|_{\ell_{2}} \le \varepsilon} \left| y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}(\boldsymbol{x} + \boldsymbol{\delta})) \right| - \left(\left| y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}) \right| + \varepsilon \|\boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\mathbb{I}(\boldsymbol{W} \boldsymbol{x} > 0)) \boldsymbol{\theta}\|_{\ell_{2}} \right) \right| = \alpha_{d},$$

where $\alpha_d = O\left(\frac{\log(d)}{d^{1/6}}\right)$.

One can thus write from (6.1) and (6.5) that for any $\theta \in C_{\theta}$

$$\left| \mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}) \right| \leq \alpha_d^2 + \alpha_d \frac{1}{n} \sum_{i=1}^n \max_{\|\boldsymbol{\delta}_i\|_{\ell_2} \leq \varepsilon} \left| y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}(\boldsymbol{x}_i + \boldsymbol{\delta}_i)) \right|$$

$$\leq \alpha_d^2 + \alpha_d \frac{1}{n} \sum_{i=1}^n \left| y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i) \right| + \alpha_d \|\boldsymbol{\theta}\|_{\ell_2} \|\boldsymbol{W}\| \varepsilon$$

$$\leq \alpha_d^2 + \alpha_d \frac{1}{n} \sum_{i=1}^n \left| y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i) \right| + c_1 \alpha_d,$$

where $c_1 > 0$ is an absolute constant. The second inequality follows from the fact that the ReLU function is 1-Lipschitz, and the third inequality follows from $\theta \in C_{\theta}$ as well as the fact that ||W|| is bounded.

We will now show that

(B.23)
$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \frac{1}{n} \sum_{i=1}^{n} \left| y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i) \right| = O_{d, \mathbb{P}}(d^{\frac{1}{12}}).$$

It is easy to see that proving the above relation will finish the proof as $\alpha_d = O\left(\frac{\log(d)}{d^{1/6}}\right)$.

Fix a θ such that $\|\theta\|_{\ell_2} \leq C$. Recall that (x_i, y_i) are generated i.i.d. according to the the distribution (2.1). Since the random variables $|y_i - \theta^T \sigma(Wx_i)|$ are sub-gaussian (see e.g. (D.93) in Lemma D.9), we can write

(B.24)
$$\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^{n}\left|y_{i}-\boldsymbol{\theta}^{\mathsf{T}}\sigma(\boldsymbol{W}\boldsymbol{x}_{i})\right|\geq d^{\frac{1}{12}}\right)\leq c_{2}e^{-c_{2}d^{\frac{7}{6}}},$$

for an absolute constant $c_2 > 0$.

Now, to prove (B.23), we use an ϵ -net argument. Consider a 1-net of the set $\{\theta : \|\theta\|_{\ell_2} \le C\}$. We know that such a 1-net \mathcal{S} exists with size at most $|\mathcal{S}| \le 2^{c_3d}$ where $c_3 > 0$ is an absolute constant. Let $\theta_1 \in \mathcal{S}$ be a vector in this net, and consider another vector θ_2 in the 1-neighborhood of θ_1 – i.e. $\|\theta_1 - \theta_2\|_{\ell_2} \le 1$. We can write

$$\left| \frac{1}{n} \sum_{i=1}^{n} |y_i - \boldsymbol{\theta}_1^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| - \frac{1}{n} \sum_{i=1}^{n} |y_i - \boldsymbol{\theta}_2^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| \right| \leq \frac{1}{n} \sum_{i=1}^{n} \left| (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_i) \right|$$

$$= \frac{1}{n} \left\| (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)^{\mathsf{T}} \boldsymbol{M} \right\|_{\ell_1}$$

$$\leq \frac{1}{\sqrt{n}} ||\boldsymbol{M}|| \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_{\ell_2}$$

$$\leq \frac{1}{\sqrt{n}} ||\boldsymbol{M}||,$$
(B.25)

where the matrix M is defined as $M = [\sigma(Wx_1)|\sigma(Wx_2)|\cdots|\sigma(Wx_n)]$. Now, since the random vectors $\sigma(Wx_i)$ are independently generated and sub-gaussian (see (D.93)), we can conclude that

$$(B.26) \mathbb{P}\left(\|M\| \ge c_4 \sqrt{n}\right) \le c_5 e^{-c_5 d},$$

for absolute constants c_4 , $c_5 > 0$ (recall that d, n, and N grow proportionally as per Assumption 1). As a result, from (B.25) and (B.26), we have

(B.27)
$$\mathbb{P}\left(\sup_{\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2}:\|\boldsymbol{\theta}_{1}-\boldsymbol{\theta}_{2}\|_{\ell_{2}}\leq 1}\left|\frac{1}{n}\sum_{i=1}^{n}\left|y_{i}-\boldsymbol{\theta}_{1}^{\mathsf{T}}\sigma(\boldsymbol{W}\boldsymbol{x}_{i})\right|-\frac{1}{n}\sum_{i=1}^{n}\left|y_{i}-\boldsymbol{\theta}_{2}^{\mathsf{T}}\sigma(\boldsymbol{W}\boldsymbol{x}_{i})\right|\right|\geq c_{4}\right)\leq c_{5}e^{-c_{5}d}.$$

Now, by using (B.24) and (B.27), and a union bound argument over S, we obtain:

$$\mathbb{P}\left(\sup_{\boldsymbol{\theta}:\|\boldsymbol{\theta}\|_{\ell_{2}} \leq C} \frac{1}{n} \sum_{i=1}^{n} |y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_{i})| \geq d^{\frac{1}{12}} + c_{4}\right) \leq \mathbb{P}\left(\sup_{\boldsymbol{\theta} \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^{n} |y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_{i})| \geq d^{\frac{1}{12}}\right) + c_{5} e^{-c_{5} d} \leq c_{2} e^{c_{3} d - c_{2} d^{\frac{7}{6}}} + c_{5} e^{-c_{5} d} = O(e^{-c_{5} d}).$$

The claim (B.23) now follows because $C_{\theta} \subseteq \{\theta : \|\theta\|_{\ell_2} \le C\}$.

APPENDIX C: PROOFS OF STEP 2: CONCENTRATION OF THE ADVERSARIAL EFFECTS

C.1. Proof of Proposition 6.4 Recall the high probability event $\mathcal{E}_{\boldsymbol{W}}$ given by (6.2). We also define the event $\mathcal{E}_{\boldsymbol{x}} \coloneqq \{\|\boldsymbol{x}_i\|_{\ell_2} \le \sqrt{5d}, \ \forall i \in [n]\}$. Since $\boldsymbol{x}_i \sim \mathsf{N}(0, \boldsymbol{I}_d), \ \|\boldsymbol{x}_i\|_{\ell_2}^2 \sim \chi_d^2$ is a chi-squared distribution with d degrees of freedom. Using chi-squared distribution tail bound (see e.g. [49, lemma 1]) along with a union bound over $i \in [n]$, we obtain $\mathbb{P}(\mathcal{E}_{\boldsymbol{x}}) \ge 1 - ne^{-d}$. Since d and n grow proportionally as per Assumption 1, both of the events $\mathcal{E}_{\boldsymbol{W}}$ and $\mathcal{E}_{\boldsymbol{x}}$ are high probability events, and so it suffices to prove the claim 6.7 on the event $\mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{x}}$.

To prove the proposition, we first state the following lemma which establishes a deviation bound for a fixed $\theta \in C_{\theta}$ and fixed $i \in [n]$.

Lemma C.1 For any fixed $\theta \in C_{\theta}$ and fixed $i \in [n]$, the following holds:

$$\mathbb{P}_{x_i} \left\{ |\eta_i(\boldsymbol{\theta})|^2 - \mathbb{E}[\eta_i(\boldsymbol{\theta})|^2] | \geq \gamma; \, \mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{x}} \right\} \leq \frac{c \log^6(d)}{d\gamma^2},$$

for some absolute constant c > 0.

Proof of Lemma C.1 is given in Section C.2.

Fix $\theta \in C_{\theta}$ and recall our notation $\nu_i(\theta; \gamma) := \mathbb{I}(|\eta_i(\theta)^2 - \mathbb{E}[\eta_i(\theta)^2]| > \gamma)$. Given that x_i are i.i.d, the random variables $\nu_i \in \{0, 1\}$ are also i.i.d. Bernoulli random variables. Therefore,

$$\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^{n}\nu_{i}(\boldsymbol{\theta};\gamma) \geq \frac{1}{\sqrt{\log(d\gamma^{2})}}\right) = \mathbb{P}\left(\sum_{i=1}^{n}\nu_{i}(\boldsymbol{\theta};\gamma) \geq \frac{n}{\sqrt{\log(d\gamma^{2})}}\right) \\
\leq \sum_{\ell=\frac{n}{\sqrt{\log(d\gamma^{2})}}}^{n} \binom{n}{\ell} \mathbb{E}[\nu_{1}(\boldsymbol{\theta};\gamma)]^{\ell} (1 - \mathbb{E}[\nu_{1}(\boldsymbol{\theta};\gamma)])^{n-\ell} \\
\leq \mathbb{E}[\nu_{1}(\boldsymbol{\theta};\gamma)]^{\frac{n}{\sqrt{\log(d\gamma^{2})}}} \sum_{\ell=\frac{n}{\sqrt{\log(d\gamma^{2})}}}^{n} \binom{n}{\ell} \\
\leq \mathbb{E}\left[\nu_{1}(\boldsymbol{\theta};\gamma)\right]^{\frac{n}{\sqrt{\log(d\gamma^{2})}}} \leq \left(2c\log^{6}(d)\right)^{n} e^{-n\sqrt{\log(d\gamma^{2})}}, \\
\leq 2^{n} \left(\frac{c\log^{6}(d)}{d\gamma^{2}}\right)^{\frac{n}{\sqrt{\log(d\gamma^{2})}}} \leq \left(2c\log^{6}(d)\right)^{n} e^{-n\sqrt{\log(d\gamma^{2})}}, \\
\leq 2^{n} \left(\frac{c\log^{6}(d)}{d\gamma^{2}}\right)^{\frac{n}{\sqrt{\log(d\gamma^{2})}}} \leq \left(2c\log^{6}(d)\right)^{n} e^{-n\sqrt{\log(d\gamma^{2})}}, \\
\leq 2^{n} \left(\frac{c\log^{6}(d)}{d\gamma^{2}}\right)^{\frac{n}{\sqrt{\log(d\gamma^{2})}}} \leq \left(2c\log^{6}(d)\right)^{n} e^{-n\sqrt{\log(d\gamma^{2})}}, \\
\leq 2^{n} \left(\frac{c\log^{6}(d)}{d\gamma^{2}}\right)^{\frac{n}{\log(d\gamma^{2})}} \leq \left(2c\log^{6}(d)\right)^{n} e^{-n\sqrt{\log(d\gamma^{2})}}, \\
\leq 2^{n} \left(\frac{c\log^{6}(d)}{d\gamma^{2}}\right)^{n} e^{-n\sqrt{\log(d\gamma^{2})}}, \\$$

where the last step follows from Lemma C.1 by which $\mathbb{E}[\nu_1(\boldsymbol{\theta};\gamma)] \leq c \log^6(d)/(d\gamma^2)$ on the event $\mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{x}}$.

Note that the above bound was for a fixed $\theta \in C_{\theta}$. In order to prove claim 6.7 we use an ε -net argument. We write

$$\sup_{\boldsymbol{\theta} \in C_{\boldsymbol{\theta}}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\boldsymbol{\theta}; \gamma) \leq \sup_{\|\boldsymbol{\theta}\|_{\ell_{2}} \leq C_{0}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\boldsymbol{\theta}; \gamma)$$

$$= \sup_{\|\boldsymbol{\theta}\|_{\ell_{2}} = C_{0}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\boldsymbol{\theta}; \gamma)$$

$$= \sup_{\boldsymbol{\theta} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\boldsymbol{\theta}; \frac{\gamma}{C_{0}^{2}}),$$
(C.2)

where the first step follows from definition of C_{θ} ; the second step follows from a simple scaling argument, and the third step follows from definition of $\eta_i^2(\theta)$ and $\nu_i(\theta;\gamma)$. We recall that \mathbb{S}^{d-1} denotes the unit (d-1)-dimensional sphere.

that \mathbb{S}^{d-1} denotes the unit (d-1)-dimensional sphere. Next consider a ε -net \mathcal{N} of \mathbb{S}^{d-1} for $\varepsilon = c_0 \gamma$. By [91, Lemma 5.2] we can choose the net \mathcal{N} so that $|\mathcal{N}| \leq (1 + \frac{2}{c_0 \gamma})^d$. We use the lemma below to relate the quantity $\nu_i(\boldsymbol{\theta}; \gamma)$ for $\boldsymbol{\theta} \in \mathbb{S}^{d-1}$ to a $\boldsymbol{\theta} \in \mathcal{N}$.

Lemma C.2 For $\theta \in \mathbb{S}^{d-1}$ choose $\widetilde{\theta} \in \mathbb{N}$ which approximates θ as $\|\theta - \widetilde{\theta}\|_{\ell_2} \leq c_0 \gamma$. On the event \mathcal{E}_W we have the following for all $i \in [n]$:

(C.3)
$$\nu_i(\boldsymbol{\theta}; \gamma) = 1 \Longrightarrow \nu_i(\widetilde{\boldsymbol{\theta}}; \gamma(1 - 2c_0\sqrt{\psi_{1,d}} - 2c_0C)) = 1.$$

We refer to Section C.3 for the proof of Lemma C.2.

Continuing from (C.2) and using Lemma C.2 we get

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\boldsymbol{\theta}; \gamma) \leq \sup_{\boldsymbol{\theta} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\boldsymbol{\theta}; \frac{\gamma}{C_{0}^{2}})$$

$$\leq \sup_{\widetilde{\boldsymbol{\theta}} \in \mathcal{N}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\widetilde{\boldsymbol{\theta}}; \frac{\gamma}{C_{0}^{2}} (1 - 2c_{0}\sqrt{\psi_{1,d}} - 2c_{0}C)).$$

Let $\tilde{\gamma} := \frac{\gamma}{C_0^2} (1 - 2c_0 \sqrt{\psi_{1,d}} - 2c_0 C)$. By choosing the constant c_0 small enough we have $\tilde{\gamma} \ge 0$. Using (C.1) along with union-bounding over the net \mathcal{N} we get

$$\mathbb{P}\left(\sup_{\widetilde{\boldsymbol{\theta}} \in \mathcal{N}} \frac{1}{n} \sum_{i=1}^{n} \nu_{i}(\widetilde{\boldsymbol{\theta}}; \widetilde{\gamma}) \geq \frac{1}{\sqrt{\log(d\widetilde{\gamma}^{2})}}\right) \leq \left(1 + \frac{2}{c_{0}\gamma}\right)^{d} \left(2c \log^{6}(d)\right)^{n} e^{-n\sqrt{\log(d\widetilde{\gamma}^{2})}}.$$

Since n and d grow proportionally and also $\gamma, \tilde{\gamma}$ are of same order, the above event is a high probability event if $\log(1/\gamma) = o(\sqrt{\log(d)})$ or equivalently if $\frac{1}{\gamma} = e^{o(\sqrt{\log(d)})}$. The result follows by combining the above bound with (C.4).

C.2. Proof of Lemma C.1 We decompose the step function as

$$\mathbb{I}(z > 0) = \mu_0 + \mu_1 z + \mu_* \varphi(z) ,$$

where for $G \sim N(0,1)$,

$$\mu_0 := \mathbb{E}[\mathbb{I}(G > 0)] = \frac{1}{2}, \quad \mu_1 = \mathbb{E}[G\mathbb{I}(G > 0)] = \frac{1}{\sqrt{2\pi}}, \quad \mu_*^2 := \mathbb{E}[\mathbb{I}(G > 0)] - \mu_0^2 - \mu_1^2 = \frac{1}{4} - \frac{1}{2\pi}.$$

Here, $\varphi(z)$ is the nonlinear component of the step function which is orthogonal to the constant and linear components in the following sense: $\mathbb{E}[\varphi(G)] = 0$ and $\mathbb{E}[G\varphi(G)] = 0$. We write

(C.5)
$$\mathbb{I}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0) = \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle + \mu_{*} u_{\ell i}, \quad \text{where: } u_{\ell i} \coloneqq \varphi(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle),$$

noting that $\langle w_{\ell}, x_i \rangle$ and $\langle w_k, x_i \rangle$ are jointly Gaussian with

$$\mathbb{E}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle^{2}) = \mathbb{E}(\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle^{2}) = 1, \quad \mathbb{E}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle \langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle) = \langle \boldsymbol{w}_{k}, \boldsymbol{w}_{\ell} \rangle.$$

Therefore, we have (see e.g., [15, Table 1])

$$\mathbb{E}[\mathbb{I}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0)\mathbb{I}(\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle > 0)] = \frac{\pi - \cos^{-1}(\langle \boldsymbol{w}_{k}, \boldsymbol{w}_{\ell} \rangle)}{2\pi}$$

$$= \frac{1}{4} + \frac{1}{2\pi} \langle \boldsymbol{w}_{\ell}, \boldsymbol{w}_{k} \rangle + O(\langle \boldsymbol{w}_{\ell}, \boldsymbol{w}_{k} \rangle^{3})$$

$$= \frac{1}{4} + \frac{1}{2\pi} \langle \boldsymbol{w}_{\ell}, \boldsymbol{w}_{k} \rangle + O(d^{-3/2} \log^{3}(d)).$$
(C.6)

To bound the correlation of variables $u_{\ell i}, u_{k i}$, we write

$$\mu_{\star}^{2} \mathbb{E}[u_{\ell i} u_{k i}] = \mathbb{E}\left[\left\{\mathbb{I}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0) - \frac{1}{2} - \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle\right\} \left\{\mathbb{I}(\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle > 0) - \frac{1}{2} - \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle\right\}\right]$$

$$= \mathbb{E}\left[\mathbb{I}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0)\mathbb{I}(\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle > 0)\right] + \mathbb{E}\left[\left(\frac{1}{2} + \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle\right)\left(\frac{1}{2} + \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle\right)\right]$$
(C.7)
$$- \mathbb{E}\left[\mathbb{I}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0)\left(\frac{1}{2} + \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle\right)\right] - \mathbb{E}\left[\mathbb{I}(\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle > 0)\left(\frac{1}{2} + \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle\right)\right].$$

The first term above is calculated in (C.6). For the second term, we have

(C.8)
$$\mathbb{E}\left[\left(\frac{1}{2} + \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i}\rangle\right)\left(\frac{1}{2} + \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i}\rangle\right)\right] = \frac{1}{4} + \frac{\langle \boldsymbol{w}_{\ell}, \boldsymbol{w}_{k}\rangle}{2\pi}.$$

For the third term we write

$$\mathbb{E}\left[\mathbb{I}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0)\left(\frac{1}{2} + \frac{1}{\sqrt{2\pi}}\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle\right)\right] = \frac{1}{4} + \frac{1}{\sqrt{2\pi}}\mathbb{E}\left[\mathbb{I}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0)\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle\right] \\
= \frac{1}{4} + \frac{1}{\sqrt{2\pi}}\mathbb{E}\left[\langle \boldsymbol{w}_{k}, \boldsymbol{x}_{i} \rangle \middle| \langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0\right]\mathbb{P}(\langle \boldsymbol{w}_{\ell}, \boldsymbol{x}_{i} \rangle > 0) \\
\stackrel{(a)}{=} \frac{1}{4} + \frac{1}{2\sqrt{2\pi}}\langle \boldsymbol{w}_{\ell}, \boldsymbol{w}_{k} \rangle \frac{\phi(0)}{1 - \Phi(0)} \\
= \frac{1}{4} + \frac{1}{2\pi}\langle \boldsymbol{w}_{\ell}, \boldsymbol{w}_{k} \rangle.$$
(C.9)

Here (a) follows from lemma below.

Lemma C.3 For Z_1 , $Z_2 \sim N(0,1)$ with $\mathbb{E}[Z_1 Z_2] = \rho$ we have

$$\mathbb{E}[Z_1|Z_2>z] = \rho \frac{\phi(z)}{(1-\Phi(z))},$$

where $\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$ is the density of standard normal and $\Phi(z) = \int_{-\infty}^{z} \phi(t) dt$ is its CDF.

Using Equations (C.6), (C.8) and (C.9) in (C.7) we obtain

(C.10)
$$\mathbb{E}[u_{\ell i}u_{ki}] = O\left(d^{-3/2}\log^3(d)\right).$$

Substituting for the sign function $\mathbb{I}(\langle w_{\ell}, x_i \rangle > 0)$ from (C.5) we get

$$\boldsymbol{W}^\mathsf{T} \mathrm{diag}(\mathbb{I}(\boldsymbol{W} \boldsymbol{x}_i > 0)) \, \boldsymbol{\theta} = \boldsymbol{W}^\mathsf{T} \mathrm{diag}(\boldsymbol{\theta}) \mathbb{I}(\boldsymbol{W} \boldsymbol{x}_i > 0)$$

(C.11)
$$= \mathbf{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \left(\frac{1}{2} \mathbf{1} + \mu_* \mathbf{u}_i + \frac{1}{\sqrt{2\pi}} \mathbf{W} \mathbf{x}_i \right),$$

with $u_i = (u_{\ell i})_{\ell=1}^N$. To lighten the notation, we use the shorthand $h_i := W^T \operatorname{diag}(\theta)(\frac{1}{2}\mathbf{1} + \mu_* u_i)$. We next decompose $\eta_i(\theta)^2$ into three terms as follows:

(C.12)
$$\eta_i(\boldsymbol{\theta})^2 = \frac{1}{2\pi} \| \boldsymbol{W}^\mathsf{T} \mathrm{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_i \|_{\ell_2}^2 + \| \boldsymbol{h}_i \|_{\ell_2}^2 + \sqrt{\frac{2}{\pi}} \langle \boldsymbol{h}_i, \boldsymbol{W}^\mathsf{T} \mathrm{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_i \rangle.$$

We next provide deviation bounds for each of these terms by putting which together we obtain the desired claim.

We start by the first term in (C.12). Note that since $\theta \in C_{\theta}$, we have the following bounds conditioned on the event \mathcal{E}_{W} :

$$\left\| \boldsymbol{W}^{\mathsf{T}} \mathrm{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{W}^{\mathsf{T}} \mathrm{diag}(\boldsymbol{\theta}) \boldsymbol{W} \right\| \leq \left\| \boldsymbol{W} \right\|^{4} \left\| \boldsymbol{\theta} \right\|_{\ell_{\infty}}^{2} = O\left(d^{-1} \log(d)\right),$$

$$\|\boldsymbol{W}^{\mathsf{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\|_{E} \leq \sqrt{\min(d,N)}\|\boldsymbol{W}^{\mathsf{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\| = O(d^{-0.5}).$$

Therefore, by applying Hanson-Wright's inequality [72] we get

(C.13)

$$\mathbb{P}\left\{\left|\left\|\boldsymbol{W}^{\mathsf{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\boldsymbol{x}_{i}\right\|_{\ell_{2}}^{2}-\mathbb{E}\left[\left\|\boldsymbol{W}^{\mathsf{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\boldsymbol{x}_{i}\right\|_{\ell_{2}}^{2}\right]\right| > \gamma \log(d); \mathcal{E}_{\boldsymbol{W}}\right\} \leq 2e^{-c\gamma^{2}d},$$

for an absolute constant c > 0.

For the second term we bound variation in the vector h_i itself from which we obtain a deviation bound on its norm $||h_i||_{\ell_2}$. We write

$$\mathbb{E}\left[\left\|\boldsymbol{h}_{i}-\mathbb{E}[\boldsymbol{h}_{i}]\right\|_{\ell_{2}}^{2}\right]=\mu_{*}^{2}\mathbb{E}\left[\left\|\boldsymbol{W}^{\mathsf{T}}\mathrm{diag}(\boldsymbol{\theta})\boldsymbol{u}_{i}\right\|_{\ell_{2}}^{2}\right]$$

$$(C.14) \qquad = \sum_{\ell,k} \langle \boldsymbol{w}_{\ell}, \boldsymbol{w}_{k} \rangle \theta_{\ell} \theta_{k} \mathbb{E}[u_{\ell i} u_{k i}] \leq \sum_{i,j} C \frac{1}{\sqrt{d}} \|\boldsymbol{\theta}\|_{\ell_{\infty}}^{2} d^{-1.5} \log^{3}(d) = O\left(d^{-1} \log^{4}(d)\right),$$

where we used the assumption $\theta \in \mathcal{C}_{\theta}$ along with (C.10). We write $h_i = \mathbb{E}[h_i] + \delta$ and define the event $\mathcal{E}_{\delta} := \{\|\delta\|_{\ell_2} \le \gamma\}$. Therefore by using Markov's inequality along with (C.14) we obtain $\mathbb{P}(\mathcal{E}_{\delta}) \ge 1 - c \frac{\log^4(d)}{d\gamma^2}$. Furthermore,

(C.15)
$$\|\boldsymbol{h}_i\|_{\ell_2}^2 = \|\mathbb{E}[\boldsymbol{h}_i]\|_{\ell_2}^2 + \|\boldsymbol{\delta}\|_{\ell_2}^2 + 2\langle \boldsymbol{\delta}, \mathbb{E}[\boldsymbol{h}_i] \rangle.$$

On the event $\mathcal{E}_{\boldsymbol{W}}$, we have $\|\mathbb{E}[\boldsymbol{h}_i]\|_{\ell_2} = \frac{1}{2} \|\boldsymbol{W}^\mathsf{T} \operatorname{diag}(\boldsymbol{\theta}) \mathbf{1}\|_{\ell_2} \leq \frac{1}{2} \|\boldsymbol{W}\| \|\boldsymbol{\theta}\|_{\infty} \sqrt{N} \leq C \sqrt{\log(d)}$, and so $|\langle \boldsymbol{\delta}, \mathbb{E}[\boldsymbol{h}_i] \rangle| \leq C \|\boldsymbol{\delta}\|_{\ell_2}$. Hence, on the event $\mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{\delta}}$,

$$\left| \|\boldsymbol{h}_i\|_{\ell_2}^2 - \|\mathbb{E}[\boldsymbol{h}_i]\|_{\ell_2}^2 \right| \leq \gamma^2 + 2C\sqrt{\log(d)}\gamma = O\left(\gamma\sqrt{\log(d)}\right).$$

This implies that $\|\mathbb{E}[\boldsymbol{h}_i]\|_{\ell_2}^2 = \mathbb{E}[\|\boldsymbol{h}_i\|_{\ell_2}^2] + O(\gamma \log(d))$, and therefore

(C.16)
$$\left\| \left\| \boldsymbol{h}_{i} \right\|_{\ell_{2}}^{2} - \mathbb{E}\left[\left\| \boldsymbol{h}_{i} \right\|_{\ell_{2}}^{2} \right] \right\| = O\left(\gamma \sqrt{\log(d)} \right).$$

We next proceed to the third term in (C.12).

$$\langle \boldsymbol{h}_{i}, \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_{i} \rangle = \langle \mathbb{E}[\boldsymbol{h}_{i}], \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_{i} \rangle + \langle \boldsymbol{\delta}, \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_{i} \rangle$$

$$= \frac{1}{2} \langle \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{1}, \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_{i} \rangle + \langle \boldsymbol{\delta}, \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_{i} \rangle.$$
(C.17)

Note that on the event \mathcal{E}_{W} the first term above is a Lipschitz continuous function of the Gaussian vector x_i with Lipschitz constant

$$L = \|\mathbf{1}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W}\|_{\ell_{0}} \leq \sqrt{N} \|\boldsymbol{\theta}\|_{\ell_{\infty}}^{2} \|\boldsymbol{W}\|^{3} = O\left(\log(d)/\sqrt{d}\right).$$

By Gaussian isoperimetry [50], we have

(C.18)

$$\mathbb{P}\left(\left|\langle \boldsymbol{W}^{\mathcal{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{1},\boldsymbol{W}^{\mathcal{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\boldsymbol{x}_{i}\rangle-\mathbb{E}[\langle \boldsymbol{W}^{\mathcal{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{1},\boldsymbol{W}^{\mathcal{T}}\operatorname{diag}(\boldsymbol{\theta})\boldsymbol{W}\boldsymbol{x}_{i}\rangle]\right|\geq\gamma\log(d);\mathcal{E}_{\boldsymbol{W}}\right)\leq2e^{-c\gamma^{2}d},$$

for some constant c > 0. For the second term of (C.17), note that on the event $\mathcal{E}_{\delta} \cap \mathcal{E}_{W} \cap \mathcal{E}_{x}$,

(C.19)

$$|\langle \boldsymbol{\delta}, \boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_i \rangle| \leq \|\boldsymbol{\delta}\|_{\ell_2} \|\boldsymbol{W}^{\mathsf{T}} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_i\|_{\ell_2} \leq \|\boldsymbol{\delta}\|_{\ell_2} \|\boldsymbol{W}\|^2 \|\boldsymbol{\theta}\|_{\ell_{\infty}} \|\boldsymbol{x}_i\|_{\ell_2} = O(\gamma \log(d)).$$

Combining (C.18) and (C.19) with the decomposition (C.17) we get that on the event $\mathcal{E}_{\delta} \cap \mathcal{E}_{W} \cap \mathcal{E}_{x}$,

(C.20)
$$\left| \langle \boldsymbol{h}_i, \boldsymbol{W}^\mathsf{T} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_i \rangle - \mathbb{E} \left[\langle \boldsymbol{h}_i, \boldsymbol{W}^\mathsf{T} \operatorname{diag}(\boldsymbol{\theta}) \boldsymbol{W} \boldsymbol{x}_i \rangle \right] \right| = O(\gamma \log(d)),$$

with probability at least $1 - 2e^{-c\gamma^2 d}$. Putting together the deviation bounds for the three terms, given by (C.13), (C.16) and (C.20), we arrive at

$$\mathbb{P}\left(|\eta_i(\boldsymbol{\theta})^2 - \mathbb{E}[\eta_i(\boldsymbol{\theta})^2]| > C\gamma \log(d); \; \mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{x}}\right) \leq 4e^{-c\gamma^2 d} + \frac{c \log^4(d)}{d\gamma^2} = O\left(\frac{\log^4(d)}{d\gamma^2}\right).$$

Note that the above relation holds for any $\gamma > 0$. The result of the lemma now follows by letting $\gamma \leftarrow C\gamma \log(d)$.

C.3. Proof of Lemma C.2 Define the matrix

$$A := \operatorname{diag}(\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0))\boldsymbol{W}\boldsymbol{W}^{\mathsf{T}}\operatorname{diag}(\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0)).$$

By triangle inequality we have

$$\begin{aligned} |\eta_{i}(\boldsymbol{\theta})^{2} - \eta_{i}(\widetilde{\boldsymbol{\theta}})^{2}| &= |\langle \boldsymbol{A}\boldsymbol{\theta}, \boldsymbol{\theta} \rangle - \langle \boldsymbol{A}\widetilde{\boldsymbol{\theta}}, \widetilde{\boldsymbol{\theta}} \rangle| \\ &= |\langle \boldsymbol{A}\boldsymbol{\theta}, \boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}} \rangle + \langle \boldsymbol{A}(\boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}}), \widetilde{\boldsymbol{\theta}} \rangle| \\ &\leq ||\boldsymbol{A}|| \, ||\boldsymbol{\theta}||_{\ell_{2}} \, ||\boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}}||_{\ell_{2}} + ||\boldsymbol{A}|| \, ||\boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}}||_{\ell_{2}} \, ||\widetilde{\boldsymbol{\theta}}||_{\ell_{2}} \leq 2c_{0}\gamma \, ||\boldsymbol{A}|| , \end{aligned}$$

where in the last step we used the fact that $\boldsymbol{\theta}, \widetilde{\boldsymbol{\theta}} \in \mathcal{S}^{d-1}$ and $\|\boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}}\|_{\ell_2} \le c_0 \gamma$. So it suffices to bound $\|\boldsymbol{A}\|$. By definition, the matrix \boldsymbol{A} is obtained by selecting a subset of rows and columns of $\boldsymbol{W}\boldsymbol{W}^\mathsf{T}$ and replacing them with zeros. Therefore, $\|\boldsymbol{A}\| \le \|\boldsymbol{W}\boldsymbol{W}^\mathsf{T}\| \le \sqrt{\psi_{1,d}} + C$, on the event $\mathcal{E}_{\boldsymbol{W}}$.

Since the above bound in Lemma C.2 holds for any vector x_i , a similar bound also holds if the terms are replaced by their expectation with respect to x_i , whence we obtain $|\mathbb{E}[\eta_i(\boldsymbol{\theta})^2] - \mathbb{E}[\eta_i(\widetilde{\boldsymbol{\theta}})^2]| \leq 2c_0\gamma(\sqrt{\psi_{1,d}} + C)$.

By definition of $\nu_i(\boldsymbol{\theta}; \gamma)$ we have

$$\nu_{i}(\boldsymbol{\theta}; \gamma) = 1 \Longrightarrow |\eta_{i}(\boldsymbol{\theta})^{2} - \mathbb{E}[\eta_{i}(\boldsymbol{\theta})^{2}]| \geq \gamma$$

$$\Longrightarrow |\eta_{i}(\widetilde{\boldsymbol{\theta}})^{2} - \mathbb{E}[\eta_{i}(\widetilde{\boldsymbol{\theta}})^{2}]| + \left| (\eta_{i}(\boldsymbol{\theta})^{2} - \mathbb{E}[\eta_{i}(\boldsymbol{\theta})^{2}]) - (\eta_{i}(\widetilde{\boldsymbol{\theta}})^{2} - \mathbb{E}[\eta_{i}(\widetilde{\boldsymbol{\theta}})^{2}]) \right| \geq \gamma$$

$$\Longrightarrow |\eta_{i}(\widetilde{\boldsymbol{\theta}})^{2} - \mathbb{E}[\eta_{i}(\widetilde{\boldsymbol{\theta}})^{2}]| \geq \gamma - 2c_{0}\gamma(\sqrt{\psi_{1,d}} + C)$$
(C.21)
$$\Longrightarrow \nu_{i}(\widetilde{\boldsymbol{\theta}}; \gamma(1 - 2c_{0}\sqrt{\psi_{1,d}} - 2c_{0}C)) = 1.$$

C.3.1. Proof of Lemma C.3 The conditional distribution of Z_1 given Z_2 is

$$Z_1|Z_2 = z_2 \sim N(\rho z_2, (1-\rho^2)).$$

Therefore, $\mathbb{E}[Z_1|Z_2=z_2]=\rho z_2$ and so

$$\mathbb{E}[Z_1|Z_2 > z] = \mathbb{E}[Z_1|Z_2 = z_2] \mathbb{P}(Z_2 = z_2|Z_2 > z) dz_2 = \rho \mathbb{E}[Z_2|Z_2 > z].$$

Using the properties of the expectation of a truncated normal distribution, we have

$$\mathbb{E}[Z_2|Z_2>z] = \frac{\phi(z)}{(1-\Phi(z))},$$

which completes the proof.

C.4. Proof of Lemma 6.6 By (6.6) it suffices to show that

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}_{\boldsymbol{\theta}}} \frac{|\overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}) - \overset{\circ}{\mathcal{L}}(\boldsymbol{\theta})|}{1 + \min(\overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}), \overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}))} = o_{d,\mathbb{P}}(1).$$

To lighten the notation we define the shorthand $\alpha_i := |y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{W}^\mathsf{T} \mathrm{diag}\mathbb{I}(\boldsymbol{W} \boldsymbol{x}_i > 0) \boldsymbol{\theta}\|_{\ell_2}$ and so $\overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}) = 1/(2n) \sum_{i=1}^n \alpha_i^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta}$. We write $\overset{\circ\circ}{\mathcal{L}}(\boldsymbol{\theta}) = 1/(2n) \sum_{i=1}^n (\alpha_i + \beta_i)^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta}$ with

$$\beta_i := \varepsilon \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_2} - \varepsilon \| \boldsymbol{W}^\mathsf{T} \operatorname{diagl}(\boldsymbol{W} \boldsymbol{x}_i > 0) \boldsymbol{\theta} \|_{\ell_2}.$$

Since for any two positive values a, b we have $|a - b| \le \sqrt{|a^2 - b^2|}$, we can write

(C.22)
$$|\beta_i| \leq \varepsilon \Big(\| \boldsymbol{W}^\mathsf{T} \operatorname{diag}\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i > 0) \boldsymbol{\theta} \|_{\ell_2}^2 - \| \boldsymbol{J}\boldsymbol{\theta} \|_{\ell_2}^2 \Big)^{1/2} = [\eta_i(\boldsymbol{\theta})^2 - \mathbb{E}[\eta_i(\boldsymbol{\theta})^2]]^{1/2}.$$

Note that on the event $\mathcal{E}_{\boldsymbol{W}}$, $\|\boldsymbol{W}\|$ is bounded and so $\|\boldsymbol{W}^{\mathsf{T}}\mathrm{diag}\mathbb{I}(\boldsymbol{W}\boldsymbol{x}_i>0)\|$ is also bounded. Since $\boldsymbol{\theta}\in\mathcal{C}_{\boldsymbol{\theta}}$, we have $\|\boldsymbol{\theta}\|_{\ell_2}=O(1)$, which along with Lemma F.3 imply that $\max_{i\in[n]}|\eta_i(\boldsymbol{\theta})|$ and $\max_{i\in[n]}|\mathbb{E}[\eta_i(\boldsymbol{\theta})]|$ are both $O_{d,\mathbb{P}}(1)$. Therefore: (i) defining $\boldsymbol{\beta}=(\beta_1,\ldots,\beta_n)$, we have $\|\boldsymbol{\beta}\|_{\ell_2}=O_{d,\mathbb{P}}(1)$; (ii) Using (C.22) along with Corollary 6.5, we get $\frac{1}{n}|\{i:|\beta_i|>\frac{1}{\sqrt{\log(d)}}\}|=O_{d,\mathbb{P}}(1)$.

From the above we can deduce that $\frac{1}{n} \|\beta\|_{\ell_2}^2 = o_{d,\mathbb{P}}(1)$. We also define $\alpha = (\alpha_1, \dots, \alpha_n)$. For any $\theta \in \mathcal{C}_{\theta}$ we have

$$|\hat{\mathcal{L}}(\boldsymbol{\theta}) - \hat{\mathcal{L}}(\boldsymbol{\theta})| = \left| \frac{1}{2n} \sum_{i=1}^{n} (\alpha_{i} + \beta_{i})^{2} - \sum_{i=1}^{n} \alpha_{i}^{2} \right|$$

$$= \frac{\|\boldsymbol{\beta}\|_{\ell_{2}}^{2}}{2n} + \frac{1}{n} |\sum_{i=1}^{n} \alpha_{i} \beta_{i}|$$

$$\leq \frac{\|\boldsymbol{\beta}\|_{\ell_{2}}^{2}}{2n} + \frac{1}{n} \|\boldsymbol{\beta}\|_{\ell_{2}} \|\boldsymbol{\alpha}\|_{\ell_{2}}$$

$$\leq \frac{\|\boldsymbol{\beta}\|_{\ell_{2}}^{2}}{2n} + \frac{\|\boldsymbol{\beta}\|_{\ell_{2}}}{\sqrt{n}} \frac{\|\boldsymbol{\alpha}\|_{\ell_{2}}}{\sqrt{n}}$$

$$\leq \frac{\|\boldsymbol{\beta}\|_{\ell_{2}}^{2}}{2n} + \frac{\|\boldsymbol{\beta}\|_{\ell_{2}}}{2\sqrt{n}} \left(1 + \frac{\|\boldsymbol{\alpha}\|_{\ell_{2}}^{2}}{n}\right)$$

$$\leq \frac{\|\boldsymbol{\beta}\|_{\ell_2}^2}{2n} + \frac{\|\boldsymbol{\beta}\|_{\ell_2}}{2\sqrt{n}} + \frac{\|\boldsymbol{\beta}\|_{\ell_2}}{\sqrt{n}} \mathring{\mathcal{L}}(\boldsymbol{\theta})$$
$$= o_{d,\mathbb{P}}(1)(1 + \mathring{\mathcal{L}}(\boldsymbol{\theta})),$$

where the last step holds because $\frac{\|\beta\|_{\ell_2}}{\sqrt{n}} = o_{d,\mathbb{P}}(1)$.

By a similar argument, we also get

$$|\overset{\circ}{\mathcal{L}}(\boldsymbol{\theta}) - \overset{\circ}{\mathcal{L}}(\boldsymbol{\theta})| \leq o_{d,\mathbb{P}}(1)(1 + \overset{\circ\circ}{\mathcal{L}}(\boldsymbol{\theta})).$$

Combining these two bounds we get $|\overset{\circ}{\mathcal{L}}(\theta) - \overset{\circ}{\mathcal{L}}(\theta)| \le o_{d,\mathbb{P}}(1) \left(1 + \min(\overset{\circ}{\mathcal{L}}(\theta),\overset{\circ}{\mathcal{L}}(\theta))\right)$.

We next proceed to the second part. By optimality of $\widehat{\theta}$ and $\widehat{\theta}$ we have

$$(\mathbf{C}.23) \qquad \overset{\circ\circ}{\mathcal{L}}(\widehat{\boldsymbol{\theta}}) < \overset{\circ\circ}{\mathcal{L}}(\mathbf{0}) = \frac{1}{n} \sum_{i=1}^{n} y_i^2 = O_{d,\mathbb{P}}(1), \quad \mathcal{L}(\widehat{\boldsymbol{\theta}}) < \mathcal{L}(\mathbf{0}) = \frac{1}{n} \sum_{i=1}^{n} y_i^2 = O_{d,\mathbb{P}}(1).$$

As shown in Proposition 6.2, $\widehat{\theta} \in C_{\theta}$ with high probability. Likewise we have $\widehat{\theta}^* \in C_{\theta}$, with high probability (this follows from Lemma D.7 for the special case of k = n in that lemma.) Therefore using the first part of the current lemma,

$$|\overset{\circ\circ}{\mathcal{L}}(\widehat{\boldsymbol{\theta}}) - \mathcal{L}(\widehat{\boldsymbol{\theta}})| = o_{d,\mathbb{P}}(1), \quad |\overset{\circ\circ}{\mathcal{L}}(\widehat{\boldsymbol{\theta}}^*) - \mathcal{L}(\widehat{\boldsymbol{\theta}}^*)| = o_{d,\mathbb{P}}(1).$$

We therefore obtain

$$0 \leq \mathcal{L}(\widehat{\boldsymbol{\theta}}^*) - \mathcal{L}(\widehat{\boldsymbol{\theta}}) < (\mathcal{L}(\widehat{\boldsymbol{\theta}}^*) - \mathcal{L}(\widehat{\boldsymbol{\theta}}^*)) + \underbrace{(\mathcal{L}(\widehat{\boldsymbol{\theta}}^*) - \mathcal{L}(\widehat{\boldsymbol{\theta}}))}_{\leq 0} + (\mathcal{L}(\widehat{\boldsymbol{\theta}}) - \mathcal{L}(\widehat{\boldsymbol{\theta}})) \leq o_{d,\mathbb{P}}(1).$$

Since $\mathcal{L}(\boldsymbol{\theta})$ is $\frac{\zeta}{2}$ -strongly convex we have

$$\|\widehat{\boldsymbol{\theta}}^* - \widehat{\boldsymbol{\theta}}\|_{\ell_2} \le o_{d,\mathbb{P}}(1)/\zeta \to 0$$
, as $d \to \infty$.

C.5. Proof of Lemma 6.7 We define

$$\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta}) \coloneqq \mathbb{E} \left[\left(|y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x})| + \varepsilon_{\text{test}} \left\| \boldsymbol{W}^\mathsf{T} \text{diagl}(\boldsymbol{W} \boldsymbol{x} > 0) \; \boldsymbol{\theta} \right\|_{\ell_2} \right)^2 \right].$$

As an immediate result of Proposition 6.1, we have $\sup_{\theta \in C_{\theta}} |\overset{\circ}{\mathsf{AR}}(\theta) - \mathsf{AR}(\theta)| = o_d(1)$. Therefore, it suffices to show that

(C.24)
$$\sup_{\boldsymbol{\theta} \in C_{\boldsymbol{\theta}}} \frac{|\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta}) - \overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta})|}{\sqrt{\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta})}} = o_{d,\mathbb{P}}(1).$$

By expanding the terms in $\mathring{\mathsf{AR}}(\boldsymbol{\theta})$ and invoking our notation $\eta_i(\boldsymbol{\theta}) = \|\boldsymbol{W}^\mathsf{T} \mathrm{diag}\mathbb{I}(\boldsymbol{W}\boldsymbol{x} > 0) \, \boldsymbol{\theta}\|_{\ell_2}$, we have

$$\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta}) = \mathbb{E}[(y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}))^{2}] + \varepsilon_{\text{test}}^{2} \mathbb{E}[\eta_{i}(\boldsymbol{\theta})^{2}] + 2\varepsilon_{\text{test}} \mathbb{E}[|y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x})| \eta_{i}(\boldsymbol{\theta})].$$

Likewise we have

$$\overset{\circ\circ}{\mathsf{AR}}(\boldsymbol{\theta}) = \mathbb{E}[(y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}))^2] + \varepsilon_{\text{test}}^2 \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2}^2 + 2\varepsilon_{\text{test}} \mathbb{E}[|y - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x})|] \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2}.$$

Recall that by definition of J we have $\mathbb{E}[\eta_i(\boldsymbol{\theta})^2] = \|J\boldsymbol{\theta}\|_{\ell_2}^2$. Hence,

$$\left| \overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta}) - \overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta}) \right| = 2\varepsilon_{\text{test}} \left| \mathbb{E} \left[|y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x})| \left(\eta_{i}(\boldsymbol{\theta}) - \sqrt{\mathbb{E}[\eta_{i}(\boldsymbol{\theta})^{2}]} \right) \right] \right|$$

$$\leq 2\varepsilon_{\text{test}} \mathbb{E} \left[\left(y - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}) \right)^{2} \right]^{1/2} \mathbb{E} \left[\left(\eta_{i}(\boldsymbol{\theta}) - \sqrt{\mathbb{E}[\eta_{i}(\boldsymbol{\theta})^{2}]} \right)^{2} \right]^{1/2}$$

$$\leq 2\varepsilon_{\text{test}} \sqrt{\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta})} \mathbb{E} \left[\left(\eta_{i}(\boldsymbol{\theta}) - \sqrt{\mathbb{E}[\eta_{i}(\boldsymbol{\theta})^{2}]} \right)^{2} \right]^{1/2} .$$

$$(C.25)$$

To bound the right-hand side, note that for any two positive values a, b we have $(a - b)^2 \le |a^2 - b^2|$. Therefore,

(C.26)
$$\mathbb{E}\left[\left(\eta_i(\boldsymbol{\theta}) - \sqrt{\mathbb{E}[\eta_i(\boldsymbol{\theta})^2]}\right)^2\right] \le \mathbb{E}\left[\left|\eta_i(\boldsymbol{\theta})^2 - \mathbb{E}[\eta_i(\boldsymbol{\theta})^2]\right|\right].$$

Also recall that for any non-negative random variable Z, we have $\mathbb{E}[Z] = \int_0^\infty \mathbb{P}(Z \ge z) dz$. Therefore,

$$\mathbb{E}\left[\left|\eta_{i}(\boldsymbol{\theta})^{2} - \mathbb{E}\left[\eta_{i}(\boldsymbol{\theta})^{2}\right]\right|\right] = \mathbb{E}\left[\left|\eta_{i}(\boldsymbol{\theta})^{2} - \mathbb{E}\left[\eta_{i}(\boldsymbol{\theta})^{2}\right]\right|; \mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{x}}\right] + \mathbb{P}\left(\left(\mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{x}}\right)^{c}\right) \\
\int_{0}^{\infty} \mathbb{P}\left(\left|\eta_{i}(\boldsymbol{\theta})^{2} - \mathbb{E}\left[\eta_{i}(\boldsymbol{\theta})^{2}\right]\right| \geq \gamma\right) d\gamma + \mathbb{P}\left(\left(\mathcal{E}_{\boldsymbol{W}} \cap \mathcal{E}_{\boldsymbol{x}}\right)^{c}\right) \\
\leq \int_{0}^{\infty} \min\left(\frac{c}{d\gamma^{2}}, 1\right) d\gamma + c \exp\left(-\log^{2}(d)/c\right) + ne^{-d}$$
(C.27)

where the inequality follows from Lemma C.1. Next, we have

$$\int_{0}^{\infty} \min\left(\frac{c\log^{6}(d)}{d\gamma^{2}}, 1\right) d\gamma = \int_{0}^{\sqrt{c\log^{6}(d)/d}} d\gamma + \int_{\sqrt{c\log^{6}(d)/d}}^{\infty} \frac{c\log^{6}(d)}{d\gamma^{2}} d\gamma$$
(C.28)
$$= \sqrt{\frac{c\log^{6}(d)}{d}} + \frac{c\log^{6}(d)}{d} \sqrt{\frac{d}{c\log^{6}(d)}} = 2\sqrt{\frac{c\log^{6}(d)}{d}}.$$

Combining Eqs. (C.26), (C.27) and (C.28) we arrive at

$$\mathbb{E}\left[\left(\eta_{i}(\boldsymbol{\theta}) - \sqrt{\mathbb{E}[\eta_{i}(\boldsymbol{\theta})^{2}]}\right)^{2}\right] \leq 2\sqrt{\frac{c\log^{6}(d)}{d}} + c\exp(-\log^{2}(d)/c) + ne^{-d} = o_{d}(\log^{3}(d)d^{-1/2}).$$

Using the above bound in (C.25) we get that uniformly over $\theta \in C_{\theta}$,

$$|\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta}) - \overset{\circ\circ}{\mathsf{AR}}(\boldsymbol{\theta})| \leq \sqrt{\overset{\circ}{\mathsf{AR}}(\boldsymbol{\theta})} \, o_{d,\mathbb{P}}(1) \, .$$

This completes the proof of claim (C.24).

APPENDIX D: PROOFS OF STEP 3: THE GAUSSIAN EQUIVALENCE PROPERTY

D.1. Proof of Proposition 6.8 As proved in [30, Theorem 2], under the assumptions of Proposition 6.8, $(\boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W}\boldsymbol{x}), \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x}_0)$ converges in distribution to $(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f}, \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x})$. We first show that $\mathsf{AR}(\boldsymbol{\theta}) = \mathsf{AR}_{\mathrm{nl}}(\boldsymbol{\theta}) + o_d(1)$. Recalling the definition of $\mathsf{AR}_{\mathrm{nl}}(\boldsymbol{\theta})$ given by (6.18), and plugging for $y = \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x} + \boldsymbol{\xi}$, we write

$$\overset{\circ}{\mathsf{AR}_{\mathrm{nl}}}(\boldsymbol{\theta}) = \mathbb{E}\left[\left(|\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x} - \boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{f} + \boldsymbol{\xi}| + \varepsilon \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}\right)^{2}\right]$$
(D.1)
$$= \mathbb{E}\left[\left(\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x} - \boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{f} + \boldsymbol{\xi}\right)^{2}\right] + \varepsilon \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}} \mathbb{E}\left[\left|\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x} - \boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{f} + \boldsymbol{\xi}\right|\right] + \varepsilon^{2} \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}^{2}$$

Therefore, $\stackrel{\circ}{\mathsf{AR}}_{\mathrm{nl}}(\theta)$ can be written in terms of the first and second moment of random variable $|\beta^\mathsf{T} x - \theta^\mathsf{T} f + \xi|$ which converges in distribution to $|\beta^\mathsf{T} x - \theta^\mathsf{T} \sigma(Wx) + \xi|$. To show that $\stackrel{\circ}{\mathsf{AR}}_{\mathrm{nl}}(\theta) - \stackrel{\circ}{\mathsf{AR}}(\theta) \to 0$ as $d \to \infty$, we need to show that the first and second moments of $|\beta^\mathsf{T} x - \theta^\mathsf{T} f + \xi|$ converge respectively to the first and second moments of $|\beta^\mathsf{T} x - \theta^\mathsf{T} \sigma(Wx) + \xi|$. As an application of [7, Corollary of Theorem 25.12], it suffices to show that $|\beta^\mathsf{T} x - \theta^\mathsf{T} f + \xi|$ has bounded third moment. To show this, note that by Holder's inequality, $|a + b + c|^3 \le 3(|a|^3 + |b|^3 + |c|^3)$. Furthermore, $\mathbb{E}[|\xi|^3] = 2$, $\mathbb{E}[|\beta^\mathsf{T} x|^3] = 2 \|\beta\|_{\ell_2}^3 = 2$. Hence,

(D.2)
$$\mathbb{E}\left[\left|\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x} - \boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{f} + \boldsymbol{\xi}\right|^{3}\right] \leq 3\left(16 + \mathbb{E}\left[\left|\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{f}\right|^{3}\right]\right).$$

By using the Holder's inequality again we have

$$E[|\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{f}|^{3}] = \mathbb{E}\left[\left|\frac{1}{\sqrt{2\pi}}\mathbf{1}^{\mathsf{T}}\boldsymbol{\theta} + \frac{1}{2}\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{W}\boldsymbol{x} + \sqrt{\frac{1}{4} - \frac{1}{2\pi}}\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{u}\right|^{3}\right]$$

$$\leq 3\left(\frac{1}{\sqrt{2\pi^{3}}}(\mathbf{1}^{\mathsf{T}}\boldsymbol{\theta})^{3} + \frac{1}{4}\left\|\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\theta}\right\|_{\ell_{2}}^{3} + 2(\frac{1}{4} - \frac{1}{2\pi})^{3/2}\left\|\boldsymbol{\theta}\right\|_{\ell_{2}}^{3}\right).$$
(D.3)

Now note that by our assumption $\mathsf{AR}_{\mathrm{nl}}(\boldsymbol{\theta})$ is bounded, which in conjunction with characterization (6.20) implies that $\|\boldsymbol{\theta}\|_{\ell_2}$, $\mathbf{1}^\mathsf{T}\boldsymbol{\theta}$ and $\|\boldsymbol{W}^\mathsf{T}\boldsymbol{\theta} - 2\boldsymbol{\beta}\|_{\ell_2}^2$ are bounded as $d \to \infty$. This also implies that $\|\boldsymbol{W}^\mathsf{T}\boldsymbol{\theta}\|_{\ell_2}^2 \le (\|\boldsymbol{W}^\mathsf{T}\boldsymbol{\theta} - \boldsymbol{\beta}\|_{\ell_2} + \|\boldsymbol{\beta}\|_{\ell_2})^2 \le (\|\boldsymbol{W}^\mathsf{T}\boldsymbol{\theta} - \boldsymbol{\beta}\|_{\ell_2} + 1)^2$ is bounded. Putting these together, we obtain that $\mathbb{E}\left[\left|\boldsymbol{\beta}^\mathsf{T}\boldsymbol{x} - \boldsymbol{\theta}^\mathsf{T}\boldsymbol{f} + \boldsymbol{\xi}\right|^3\right]$ is bounded, which completes the argument for showing that $\mathsf{AR}_{\mathrm{nl}}(\boldsymbol{\theta}) = \mathsf{AR}(\boldsymbol{\theta}) + o_d(1)$.

We next prove the characterization (6.20). Note that

$$y - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f} = \boldsymbol{\xi} + \boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{x} - \frac{1}{2} \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{W} \boldsymbol{x} - \sqrt{\frac{1}{4} - \frac{1}{2\pi}} \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{u}$$
$$= \boldsymbol{\xi} + \langle \boldsymbol{\beta} - \frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta}, \boldsymbol{x} \rangle - \sqrt{\frac{1}{4} - \frac{1}{2\pi}} \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{u}$$
$$\sim \mathsf{N}(0, M(\boldsymbol{\theta})^{2}),$$

with

$$M(\theta)^2 = \tau^2 + \left\| \frac{1}{2} \mathbf{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \left(\frac{1}{4} - \frac{1}{2\pi} \right) \|\boldsymbol{\theta}\|_{\ell_2}^2.$$

We then write

$$\overset{\circ}{\mathsf{AR}_{\mathrm{nl}}}(\boldsymbol{\theta}) = \mathbb{E}\left[\left(|y - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f}| + \varepsilon_{\mathrm{test}} \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_{2}}\right)^{2}\right]
= \mathbb{E}\left[\left(y - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f}\right)^{2}\right] + 2\varepsilon_{\mathrm{test}} \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_{2}} \mathbb{E}\left[\left|y - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f}\right|\right] + \varepsilon_{\mathrm{test}}^{2} \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_{2}}^{2}
= M(\boldsymbol{\theta})^{2} + \varepsilon_{\mathrm{test}}^{2} \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_{2}}^{2} + 2\sqrt{\frac{2}{\pi}} \varepsilon_{\mathrm{test}} M(\boldsymbol{\theta}) \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_{2}},$$
(D.5)

using the first and second moment of folded normal distribution.

D.2. Proof of Theorem 6.9 Before stating the proof, we remark that our proof is an adaptation of the powerful machinery developed [39]; however, since our individual adversarial losses have an additional term $\varepsilon \|J\theta\|_{\ell_2}$, there are some additional details in the proofs which we provide in the following. Also, since our activation function is not odd, we will use the CLT-type result of [30] instead of the one provided in [39].

Recall that we are seeking to analyze the asymptotic values of the following quantities:

$$\Phi_A := \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} (|y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2})^2 + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \|\frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta}\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta})^2,$$

$$\Phi_{B} \coloneqq \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} (|y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f}_{i}| + \varepsilon ||\boldsymbol{J} \boldsymbol{\theta}||_{\ell_{2}})^{2} + \lambda ||\boldsymbol{\theta}||_{\ell_{2}}^{2} + \lambda_{w} ||\frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta}||_{\ell_{2}}^{2} + \lambda_{s} \frac{\log(d)}{d} (\mathbf{1}^{\mathsf{T}} \boldsymbol{\theta})^{2},$$

To simplify our notation, and without loss of generality, we absorb the value ϵ into J and, with a slight abuse of notation, consider $J \leftarrow \epsilon J$ (and hence the eigenvalues of the matrix J depend on ϵ). Hence, above the quantities of interest become

$$\Phi_{A} := \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} (|y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_{i})| + ||\boldsymbol{J} \boldsymbol{\theta}||_{\ell_{2}})^{2} + \lambda ||\boldsymbol{\theta}||_{\ell_{2}}^{2} + \lambda_{w} ||\frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta}||_{\ell_{2}}^{2} + \lambda_{s} \frac{\log(d)}{d} (\boldsymbol{1}^{\mathsf{T}} \boldsymbol{\theta})^{2},$$

$$\Phi_{B} := \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} \left(\| y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{f}_{i} \| + \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_{2}} \right)^{2} + \lambda \| \boldsymbol{\theta} \|_{\ell_{2}}^{2} + \lambda_{w} \left\| \frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_{2}}^{2} + \lambda_{s} \frac{\log(d)}{d} (\boldsymbol{1}^{\mathsf{T}} \boldsymbol{\theta})^{2},$$

For technical reasons, we first need to make the objectives smooth. We thus define $g(x) = \sqrt{x+\gamma}$, and define the smoothed loss

(D.6)
$$\ell(\boldsymbol{\theta}; \boldsymbol{r}, y) = (\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y)^2 + \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_2}^2 + 2g\left((\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y)^2 \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_2}^2\right),$$

Note that when $\gamma = 0$, we have $\ell(\boldsymbol{\theta}; \sigma(\boldsymbol{W}\boldsymbol{x}_i), y_i) = (|y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W}\boldsymbol{x}_i)| + ||\boldsymbol{J}\boldsymbol{\theta}||_{\ell_2})^2$ and $\ell(\boldsymbol{\theta}; \boldsymbol{f}_i, y_i) = (|y_i - \boldsymbol{\theta}^\mathsf{T} \boldsymbol{f}_i + ||\boldsymbol{J}\boldsymbol{\theta}||_{\ell_2})^2$.

In the following, we consider an arbitrary but *fixed* value $\gamma > 0$, and with some abuse of notation, we let

$$\Phi_A := \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\theta}; \sigma(\boldsymbol{W}\boldsymbol{x}_i), y_i) + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \left\| \frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta})^2,$$

$$\Phi_{B} \coloneqq \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} \ell(\boldsymbol{\theta}; \boldsymbol{f}_{i}, y_{i}) + \lambda \|\boldsymbol{\theta}\|_{\ell_{2}}^{2} + \lambda_{w} \left\| \frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_{2}}^{2} + \lambda_{s} \frac{\log(d)}{d} (\mathbf{1}^{\mathsf{T}} \boldsymbol{\theta})^{2},$$

For these quantities Φ_A, Φ_B , we show then show in the following that the statement of the Theorem is true. Then, the result of the Theorem for the original losses – i.e. when $\gamma = 0$ – follows simply by taking the limit $\gamma \to 0$ (note that this limit is taken *after* the limit $d \to \infty$; also, note that $\sup_{x \ge 0} \{g(x) - x\} = \sqrt{\gamma}$).

We will use the Lindeberg's leave-one-out technique. In a nutshell, we start with the quantity Φ_B , and through n consecutive steps, we reach to the quantity Φ_A . In the k-th step, we will replace the feature vector \mathbf{f}_k with $\sigma(\mathbf{W}\mathbf{x}_k)$. We will then show that each of these replacements has a negligible effect (i.e. $o_n(1)/n$) on our quantities of interest, leading to the proof of the theorem.

Let us now proceed with the details. The proof has multiple steps which will be put together in Section D.2.6 to obtain the proof of the theorem.

We begin by defining

$$\Phi_k = \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^k \ell(\boldsymbol{\theta}; \sigma(\boldsymbol{W} \boldsymbol{x}_i)_i, y_i) + \frac{1}{n} \sum_{i=k+1}^n \ell(\boldsymbol{\theta}; \boldsymbol{f}_i, y_i) + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \left\| \frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta})^2,$$

Roughly speaking, our goal is to show that for all $k \in [n]$, we have $\Phi_k \approx \Phi_{k-1} + o_n(1)/n$. To make this entirely rigorous, we need to define several new quantities and understand their relations. For $k \in [n]$ let

(D.8)
$$R_{-k}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{k-1} \ell\left(\boldsymbol{\theta}; \sigma(\boldsymbol{W}\boldsymbol{x}_i), y_i\right) + \frac{1}{n} \sum_{i=k+1}^{n} \ell\left(\boldsymbol{\theta}; \boldsymbol{f}_i, y_i\right) + \lambda \left\|\boldsymbol{\theta}\right\|_{\ell_2}^2 + \lambda_w \left\|\frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta}\right\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\mathbf{1}^{\mathsf{T}} \boldsymbol{\theta})^2,$$

and

(D.9)
$$R_k(\boldsymbol{\theta}, \boldsymbol{r}) = \frac{1}{n} \ell(\boldsymbol{\theta}; \boldsymbol{r}, y_k) + R_{-k}(\boldsymbol{\theta}).$$

Let us denote the minimizers of the above two objectives by

(D.10)
$$\theta_{-k}^* = \arg\min_{\theta} R_{-k}(\theta), \text{ and } \theta_k^*(r) = \arg\min_{\theta} R_k(\theta, r)$$

and

(D.11)
$$\Phi_{-k} = \min_{\boldsymbol{\theta}} R_{-k}(\boldsymbol{\theta}),$$

and

(D.12)
$$\Phi_k(\boldsymbol{r}) = \min_{\boldsymbol{\theta}} R_k(\boldsymbol{\theta}, \boldsymbol{r}).$$

It will also be convenient to work with approximate versions of the term $R_k(\theta, r)$ in (D.9). Hence, we define below we define $\mathbb{R}_k(\theta, r)$ which is essentially obtained by Taylor-expanding the term $R_{-k}(\theta)$ in (D.9).

(D.13)
$$S_k(\boldsymbol{\theta}, \boldsymbol{r}) = \Phi_{-k} + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \boldsymbol{H}_{-k} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) + \frac{1}{n} \ell(\boldsymbol{\theta}; \boldsymbol{r}, y_k),$$

where H_{-k} is the Hessian of $R_{-k}(\theta)$ at θ_{-k}^* , i.e.

(D.14)
$$\boldsymbol{H}_{-k} = \nabla^2 R_{-k}(\boldsymbol{\theta}) |_{\boldsymbol{\theta} = \boldsymbol{\theta}_{-k}^*}.$$

Finally, we denote the minimizer of $S(\theta, r)$ by

(D.15)
$$\tilde{\theta}_k(r) = \arg\min_{\boldsymbol{a}} S_k(\boldsymbol{\theta}, \boldsymbol{r}),$$

and

(D.16)
$$\Psi_k(\boldsymbol{r}) = \min_{\boldsymbol{\theta}} S_k(\boldsymbol{\theta}, \boldsymbol{r}).$$

Simplification of Notation. We note from (D.7) that in our analysis the feature vectors are either $\mathbf{r}_i = \mathbf{a}_i$ or $\mathbf{r}_i = \mathbf{b}_i$; i.e. we can write Φ_k as

(D.17)
$$\Phi_k = \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\theta}; \boldsymbol{r}_i, y_i) + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \left\| \frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta})^2,$$

where the feature vectors \mathbf{r}_i they are generated according to one of the following distributions

(D.18)
$$r_i = \sigma(\mathbf{W}\mathbf{x}_i)$$
 or $r_i = f_i := \mu_1 \mathbf{W}\mathbf{x}_i + \mu_2 \mathbf{u}_i$

It will be sometimes easier in our analysis to use (D.17), i.e. use r_i for both $\sigma(\mathbf{W}\mathbf{x}_i)$ and f_i , but we will keep in mind that for $i \le k$ we have $r_i = f_i$ and for i > k we have $r_i = \sigma(\mathbf{W}\mathbf{x}_i)$.

Details of the Gradient and Hessian of ℓ **.** In the following, we will need to work out the first and second derivatives of the loss function ℓ , given in (D.6), at multiple points. In order to present the derivations more compactly, let us denote

(D.19)
$$h(\boldsymbol{\theta}; \boldsymbol{r}, y) \coloneqq (\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y)^2 \| \boldsymbol{J} \boldsymbol{\theta} \|_{\ell_2}^2,$$

and provide the details for the derivatives of the loss function ℓ here. Given how the function h is defined, we can write

$$\ell(\boldsymbol{\theta}; \boldsymbol{r}, y) = (\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y)^{2} + \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}^{2} + 2g(h(\boldsymbol{\theta}; \boldsymbol{r}, y))$$

Using the notation ∇ for gradient w.r.t. θ , and ∇^2 for hessian w.r.t. θ , we can write

(D.20)
$$\nabla \ell(\boldsymbol{\theta}; \boldsymbol{r}, y) = 2(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y) \boldsymbol{r} + 2 \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} \boldsymbol{\theta} + 2 \nabla h(\boldsymbol{\theta}; \boldsymbol{r}, y) g'(h(\boldsymbol{\theta}; \boldsymbol{r}, y))$$

and

(D.21)

$$\nabla^{2}\ell(\boldsymbol{\theta};\boldsymbol{r},y) = 2\left(\boldsymbol{r}\boldsymbol{r}^{\mathsf{T}} + \boldsymbol{J}^{\mathsf{T}}\boldsymbol{J} + \nabla^{2}h(\boldsymbol{\theta};\boldsymbol{r},y)g'(h(\boldsymbol{\theta};\boldsymbol{r},y)) + \nabla h(\boldsymbol{\theta};\boldsymbol{r},y)(\nabla h(\boldsymbol{\theta};\boldsymbol{r},y))^{\mathsf{T}}g''(h(\boldsymbol{\theta};\boldsymbol{r},y))\right),$$

where

(D.22)
$$\nabla h(\boldsymbol{\theta}; \boldsymbol{r}, y) = 2\left(\left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y\right) \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}^{2} \boldsymbol{r} + \left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y\right)^{2} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J}\boldsymbol{\theta}\right),$$

and

(D.23)

$$\nabla^{2}h(\boldsymbol{\theta};\boldsymbol{r},y) = 2\left(\|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}^{2}\boldsymbol{r}\boldsymbol{r}^{\mathsf{T}} + 2(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{r} - y)\boldsymbol{r}\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J} + 2(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{r} - y)\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}\boldsymbol{r}^{\mathsf{T}} + (\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{r} - y)^{2}\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}\right).$$

D.2.1. Some properties of the minimizers in (D.10) In this section, we will analyze some of the properties of the vector $\boldsymbol{\theta}_{-k}^*$ and its relation with $\tilde{\boldsymbol{\theta}}_{-k}(\boldsymbol{r})$, for $k \in [n]$. We first show some basic properties of the vectors $\boldsymbol{\theta}_{-k}^*$ and $\boldsymbol{\theta}_k^*(\boldsymbol{r})$.

Lemma D.1 Fix $k \in [n]$. The following hold with absolute constants c, C > 0:

(a) The vector θ_{-k}^* is bounded in the ℓ_2 norm:

$$(D.24) \mathbb{P}\left(\|\boldsymbol{\theta}_{-k}^*\|_{\ell_2} \ge v + C\right) \le c \exp\left(-nv^2/c\right).$$

(b) The vector $\theta_k^*(r)$ is bounded in the ℓ_2 norm:

(D.25)
$$\mathbb{P}\left(\|\boldsymbol{\theta}_{k}^{*}(\boldsymbol{r})\|_{\ell_{2}} \ge v + C\right) \le c \exp\left(-nv^{2}/c\right).$$

(c) For an independently generated vector \mathbf{r} we have

(D.26)
$$\mathbb{P}\left(|\mathbf{r}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*}(\mathbf{r})| \ge v\right) \le c \exp(-v/c).$$

(d) We also have

(D.27)
$$\left|\mathbf{1}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*}\right| \leq C\sqrt{\frac{d}{\log(d)}},$$

with probability at least $1 - ce^{-cn}$.

The proof of this lemma is provided in Section D.2.7. We now show that the distance between the two minimizers θ_{-k}^* and $\tilde{\theta}_k(r)$ is of order $O(1/\sqrt{d})$.

Lemma D.2 Fix $k \in [n]$. Assuming that r is generated independently from θ_{-k}^* and according to one of the distributions in (D.18). Then, there exist absolute constants c, c' > 0 such that

(D.28)
$$\mathbb{P}\left(\left\|\tilde{\boldsymbol{\theta}}_{-k}(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^{*}\right\|_{\ell_{2}} \ge \frac{v}{\sqrt{d}}\right) \le c \exp\left(-v^{c'}/c\right).$$

Proof We start by noting that since $\tilde{\theta}_k(r)$ is the minimizer of (D.13):

(D.29)
$$\tilde{\boldsymbol{\theta}}_{k}(\boldsymbol{r}) = \arg\min_{\boldsymbol{\theta}} \left\{ S(\boldsymbol{\theta}) := \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}} \boldsymbol{H}_{-k} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}) + \frac{1}{n} \ell(\boldsymbol{\theta}; \boldsymbol{r}, y_{k}) \right\}$$

Observe that (i) the function S is λ -strongly convex due to the fact that $R_{-k}(\theta)$ is strongly-convex, and thus its Hessian H_{-k} is a PSD matrix with smallest eigenvalue lower-bounded by λ ; (ii) $S(\theta) \ge 0$ for any θ , as H_{-k} is a PSD matrix and ℓ is always positive-valued. As a result, we can write

(D.30)
$$\|\tilde{\boldsymbol{\theta}}_{k}(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^{*}\|_{\ell_{2}}^{2} \leq \frac{1}{\lambda} S(\boldsymbol{\theta}_{-k}^{*})$$

We can then write from (D.6) that

$$S(\boldsymbol{\theta}_{-k}^*) = \frac{1}{n} \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}, y_k)$$

$$\leq \frac{1}{n} C \max\{1, (\boldsymbol{\theta}_{-k}^*^{\mathsf{T}} \boldsymbol{r} - y)^2, \|\boldsymbol{J} \boldsymbol{\theta}_{-k}^*\|_{\ell_2}^2\},$$

where C > 0 is an absolute constant. The proof now follows from the result of Lemma D.1 and the fact that ||J|| is bounded, as well as the fact that d and n grow in proportion to each other.

Given the above lemma, we can analyze the behavior of $\Psi_k(r)$, defined in (D.16), in more detail.

Lemma D.3 Fix $k \in [n]$. We have (D.31)

$$\Psi_{k}(\boldsymbol{r}) = \Phi_{-k} + \frac{1}{n} \min_{\tau_{1}} \left\{ \frac{1}{2n} \left(\frac{\partial \tilde{\ell}(\tau_{1}, 0)}{\partial \tau_{1}} \boldsymbol{r} + \frac{\partial \tilde{\ell}(\tau_{1}, 0)}{\partial \tau_{2}} \boldsymbol{p} \right)^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \left(\frac{\partial \tilde{\ell}(\tau_{1}, 0)}{\partial \tau_{1}} \boldsymbol{r} + \frac{\partial \tilde{\ell}(\tau_{1}, 0)}{\partial \tau_{2}} \boldsymbol{p} \right) + \tilde{\ell}(\tau_{1}, 0) \right\} + e,$$

where (i) we have $p^{\mathsf{T}} = 2\theta_{-k}^{*} {}^{\mathsf{T}} J^{\mathsf{T}} J$; (ii) the function $\tilde{\ell}(\tau_1, \tau_2)$ is defined as

(D.32)
$$\tilde{\ell}(\tau_1, \tau_2) := \rho_1 + \rho_2 + \rho_3 \tau_1 + \tau_1^2 + \tau_2 + g((\rho_1 + \rho_3 \tau_1 + \tau_1^2)(\rho_2 + \tau_2)),$$

with $\rho_1 = (\theta_{-k}^* \mathsf{T} r - y_k)^2$, $\rho_2 = \|J\theta_{-k}^*\|_{\ell_2}^2$, and $\rho_3 = 2(\theta_{-k}^* \mathsf{T} r - y_k)$; and (iii) the value e satisfies

$$\mathbb{P}(|e| \ge \frac{v}{d^{\frac{3}{2}}}) \le c \exp\left(-v^{c'}/c\right),$$

for absolute constants c, c' > 0.

Furthermore, assuming that τ_1^* is the minimizer of the optimization problem in (D.31), we have

(D.33)
$$\tilde{\theta} - \theta_{-k}^* = \frac{1}{n} \left(\beta_1 H_{-k}^{-1} r + \beta_2 H_{-k}^{-1} p \right) + e,$$

where β_1, β_2 depend only on τ_1^* , as well as $\mathbf{r}^\mathsf{T} \boldsymbol{\theta}_{-k}^*$, and $\| \boldsymbol{J} \boldsymbol{\theta}_{-k}^* \|_{\ell_2}$; and

(D.34)
$$\mathbb{P}\left(\max\{d^{\frac{3}{2}} \|\boldsymbol{e}\|_{\ell_{2}}, |\beta_{1}|, |\beta_{2}|, \|\boldsymbol{p}\|_{\ell_{2}}\} \ge v\right) \le c \exp(-v^{c'}/c).$$

Proof In the following, to simplify notation, we use $\tilde{\theta}$ instead of $\tilde{\theta}_k(r)$. We have from (D.13) and (D.16) that

(D.35)
$$\Psi_k(\boldsymbol{r}) = \Phi_{-k} + \min_{\boldsymbol{\theta}} \left\{ \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \boldsymbol{H}_{-k} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) + \frac{1}{n} \ell(\boldsymbol{\theta}; \boldsymbol{r}, y_k) \right\}.$$

We now decompose the term $\frac{1}{n}\ell(\theta; r, y_k)$ (see (D.6)) according to the following set of simple relations:

$$\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - \boldsymbol{y} = \boldsymbol{\theta}_{-k}^{*} \boldsymbol{r} - \boldsymbol{y} + (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}} \boldsymbol{r},$$

$$\|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}^{2} = \|\boldsymbol{J}\boldsymbol{\theta}_{-k}^{*}\|_{\ell_{2}}^{2} + 2\boldsymbol{\theta}_{-k}^{*} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}) + (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}),$$

$$(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - \boldsymbol{y})^{2} = (\boldsymbol{\theta}_{-k}^{*} \boldsymbol{r} - \boldsymbol{y})^{2} + 2(\boldsymbol{\theta}_{-k}^{*} \boldsymbol{r} - \boldsymbol{y}) \boldsymbol{r}^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}) + (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}} \boldsymbol{r}^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}),$$

As a result, it is easy to obtain the following:

$$\frac{1}{n}\ell(\boldsymbol{\theta}; \boldsymbol{r}, y_{k}) = \frac{1}{n} \left\{ \rho_{1} + \rho_{2} + \rho_{3}(\boldsymbol{r}^{\mathsf{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})) + (\boldsymbol{r}^{\mathsf{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}))^{2} + \boldsymbol{p}^{\mathsf{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}) \right\}
+ \frac{1}{n} \|\boldsymbol{J}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})\|_{\ell_{2}}^{2}
+ \frac{1}{n} g \left((\rho_{1} + \rho_{3}(\boldsymbol{r}^{\mathsf{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})) + (\boldsymbol{r}^{\mathsf{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}))^{2})(\rho_{2} + \boldsymbol{p}^{\mathsf{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}) + \|\boldsymbol{J}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})\|_{\ell_{2}}^{2}) \right),$$

where the parameters ρ_j , j = 1, 2, 3, and the vector \boldsymbol{p} are defined in the following:

(D.36)
$$\rho_{1} = (\boldsymbol{\theta}_{-k}^{*} \mathbf{T} \boldsymbol{r} - y_{k})^{2},$$

$$\rho_{2} = \|\boldsymbol{J} \boldsymbol{\theta}_{-k}^{*}\|_{\ell_{2}}^{2},$$

$$\rho_{3} = 2(\boldsymbol{\theta}_{-k}^{*} \mathbf{T} \boldsymbol{r} - y_{k}),$$

$$\boldsymbol{p}^{\mathsf{T}} = 2\boldsymbol{\theta}_{-k}^{*} \mathbf{J}^{\mathsf{T}} \boldsymbol{J}.$$

We note that none of these quantities depend on the optimization variable θ and hence can be considered as constants w.r.t. the minimization procedure in (D.35).

It will be convenient to consider the following variables:

(D.37)
$$\tau_1 = \mathbf{r}^\mathsf{T}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*), \quad \text{and} \quad \tau_2 = \mathbf{p}^\mathsf{T}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*), \quad \text{and} \quad \tau_3 = \|\boldsymbol{J}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*)\|_{\ell_2}^2.$$

As a result, we can write

$$\frac{1}{n}\ell(\boldsymbol{\theta}; \boldsymbol{r}, y_k) = \frac{1}{n} \left\{ \rho_1 + \rho_2 + \rho_3 \tau_1 + \tau_1^2 + \tau_2 + \tau_3 + g \left((\rho_1 + \rho_3 \tau_1 + \tau_1^2) (\rho_2 + \tau_2 + \tau_3) \right) \right\}$$
(D.38)
$$:= \frac{1}{n} \tilde{\ell}(\tau_1, \tau_2, \tau_3).$$

Now, from (D.35), we can write the following equation for $\tilde{\theta}$ (as it is the minimizer):

(D.39)
$$\boldsymbol{H}_{-k}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) = -\frac{1}{n} \left(\frac{\partial \tilde{\ell}}{\partial \tau_1} \boldsymbol{r} + \frac{\partial \tilde{\ell}}{\partial \tau_2} \boldsymbol{p} + \frac{\partial \tilde{\ell}}{\partial \tau_3} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \right),$$

where the partial derivatives are evaluated at τ_1, τ_2, τ_3 when $\theta = \theta$. Consequently, we have

(D.40)
$$\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^* = -\frac{1}{n} \left(\frac{\partial \tilde{\ell}}{\partial \tau_1} \boldsymbol{H}_{-k}^{-1} \boldsymbol{r} + \frac{\partial \tilde{\ell}}{\partial \tau_2} \boldsymbol{H}_{-k}^{-1} \boldsymbol{p} + \frac{\partial \tilde{\ell}}{\partial \tau_3} \boldsymbol{H}_{-k}^{-1} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \right),$$

Using the above relation, we can derive a few useful properties. First, given how ρ_j 's are defined in (D.36), and by using Lemma D.1, and since $g(x) = \sqrt{x+\gamma}$ has uniformly bounded first and second derivative, and by using the fact that the operator norm of J is bounded, it is easy to show that for absolute constants c, c' we have

(D.41)
$$\mathbb{P}\left(\max\left\{\left|\frac{\partial \tilde{\ell}}{\partial \tau_1}\right|, \left|\frac{\partial \tilde{\ell}}{\partial \tau_2}\right|, \left|\frac{\partial \tilde{\ell}}{\partial \tau_3}\right|\right\} \ge v\right) \le c \exp(-v^{c'}/c),$$

where in the above relation the partial derivatives are evaluated at τ_1, τ_2, τ_3 , when $\theta = \tilde{\theta}$.

Second, by using Lemma D.2, and the fact that the norm of the matrix J is bounded, as well as the fact that the operator norm of H_{-k}^{-1} is upper-bounded by $1/\lambda$, we have for absolute constants c, c' that

$$(D.42) \mathbb{P}\left(\left\|\frac{1}{n}\boldsymbol{H}_{-k}^{-1}\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^{*})\right\|_{\ell_{2}} \geq \frac{v}{d^{\frac{3}{2}}}\right) \leq c\exp(-v^{c'}/c).$$

Third, from the definition of p in (D.36), and by using Lemma D.1 as well as (D.40) we obtain for absolute constants c, c' that

(D.43)
$$\mathbb{P}\left(\boldsymbol{p}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) \ge \frac{v}{d^{\frac{1}{2}}}\right) \le c \exp(-v^{c'}/c).$$

The above relation shows that the value of τ_2 , evaluated at $\theta = \tilde{\theta}$, is of order $O(d^{-\frac{1}{2}})$. Fourth, we can write using Lemma D.2 that

(D.44)
$$\mathbb{P}\left(\left\|\boldsymbol{J}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^*)\right\|_{\ell_2} \ge \frac{v}{d^{\frac{1}{2}}}\right) \le c \exp(-v^{c'}/c),$$

which essentially results in τ_3 , evaluated at $\theta = \tilde{\theta}$, to be of the order $O(d^{-1})$. Finally, we note that for each of the partial derivatives we can write

$$(D.45) \qquad \mathbb{P}\left(\left|\frac{\partial \tilde{\ell}(\tau_1, \tau_2, \tau_3)}{\partial \tau_j} - \frac{\partial \tilde{\ell}(\tau_1, 0, 0)}{\partial \tau_j}\right| \ge \frac{v}{d^{\frac{1}{2}}}\right) \le c \exp(-v^{c'}/c),$$

for j = 1, 2, 3 and for τ_j 's that are evaluated at $\theta = \tilde{\theta}$.

Using the above five properties, we can conclude that

(D.46)
$$\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^* = -\frac{1}{n} \left(\frac{\partial \tilde{\ell}(\tau_1, 0, 0)}{\partial \tau_1} \boldsymbol{H}_{-k}^{-1} \boldsymbol{r} + \frac{\partial \tilde{\ell}(\tau_1, 0, 0)}{\partial \tau_2} \boldsymbol{H}_{-k}^{-1} \boldsymbol{p} \right) + \boldsymbol{e},$$

and

(D.47)
$$\frac{1}{n}\tilde{\ell}(\tau_1, \tau_2, \tau_3) = \frac{1}{n}\tilde{\ell}(\tau_1, 0, 0) + e',$$

where τ_1, τ_2, τ_3 are computed from (D.37) at $\theta = \hat{\theta}$, and

$$\mathbb{P}\left(\max\{\|e\|_{\ell_2}, |e'|\} \ge \frac{v}{d^{\frac{3}{2}}}\right) \le c \exp(-v^{c'}/c).$$

As a result, by defining

$$\tilde{\ell}(\tau_1,\tau_2) \coloneqq \tilde{\ell}(\tau_1,\tau_2,0),$$

where $\tilde{\ell}$ is given in (D.38), and by plugging the solution (D.46) into the optimization in (D.35), we obtain

(D.48)

$$\Psi_{k}(\boldsymbol{r}) = \Phi_{-k} + \frac{1}{n} \min_{\tau_{1}} \left\{ \frac{1}{2n} \left(\frac{\partial \tilde{\ell}(\tau_{1},0)}{\partial \tau_{1}} \boldsymbol{r} + \frac{\partial \tilde{\ell}(\tau_{1},0)}{\partial \tau_{2}} \boldsymbol{p} \right)^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \left(\frac{\partial \tilde{\ell}(\tau_{1},0)}{\partial \tau_{1}} \boldsymbol{r} + \frac{\partial \tilde{\ell}(\tau_{1},0)}{\partial \tau_{2}} \boldsymbol{p} \right) + \tilde{\ell}(\tau_{1},0) \right\} + e,$$

where

$$\mathbb{P}\left(|e| \ge \frac{v}{d^{\frac{3}{2}}}\right) \le c \exp(-v^{c'}/c).$$

D.2.2. Bounding
$$\|\tilde{\boldsymbol{\theta}}(\boldsymbol{r}) - \boldsymbol{\theta}^*(\boldsymbol{r})\|_{\ell_2}$$

Lemma D.4 Fix $k \in [n]$. There exist absolute constants b, c, c' > 0 such that for any $v \ge 0$

$$(D.49) \qquad \mathbb{P}\left(\left\|\tilde{\boldsymbol{\theta}}_{k}(\boldsymbol{r}) - \boldsymbol{\theta}_{k}^{*}(\boldsymbol{r})\right\|_{\ell_{2}} \ge v \frac{\left(\log(d)\right)^{b}}{d}\right) \le c \exp\left(-\frac{v^{c'}}{c}\right) + c \exp\left(-(\log(d))^{2}/c\right)$$

Proof To simplify notation, in this lemma we use θ^* instead of $\theta_k^*(r)$ and $\tilde{\theta}$ instead of $\tilde{\theta}_k(r)$. Also, we define

$$q(\boldsymbol{\theta}) = \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \left\| \frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\boldsymbol{1}^\mathsf{T} \boldsymbol{\theta})^2$$

Due to the fact that $R(\theta, r)$ is λ -strongly convex, we can write

(D.50)
$$\|\boldsymbol{\theta}^* - \tilde{\boldsymbol{\theta}}\|_{\ell_2} \le \frac{1}{\lambda} \|\nabla_{\boldsymbol{\theta}} R(\tilde{\boldsymbol{\theta}}, \boldsymbol{r})\|_{\ell_2}$$

Consequently, we will bound the right-hand-side in the above relation. By using the fact that θ_{-k}^* is the minimizer of the function R_{-k} , we can write

(D.51)
$$\nabla_{\boldsymbol{\theta}} R(\tilde{\boldsymbol{\theta}}, \boldsymbol{r}) = \nabla_{\boldsymbol{\theta}} R(\tilde{\boldsymbol{\theta}}, \boldsymbol{r}) - \nabla_{\boldsymbol{\theta}} R_{-k}(\boldsymbol{\theta}_{-k}^*).$$

From (D.9) and (D.51) we have

$$\nabla_{\boldsymbol{\theta}} R(\tilde{\boldsymbol{\theta}}, \boldsymbol{r}) = \nabla R_{-k}(\tilde{\boldsymbol{\theta}}) + \frac{1}{n} \nabla \ell(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}, y_k) - \nabla R_{-k}(\boldsymbol{\theta}_{-k}^*)$$

$$= \frac{1}{n} \sum_{t+k} (\nabla \ell(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_t, y_t) - \nabla \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)) + \nabla q(\tilde{\boldsymbol{\theta}}) - \nabla q(\boldsymbol{\theta}_{-k}^*) + \frac{1}{n} \nabla \ell(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}, y_k),$$

By using the fact that $\tilde{\theta}$ is the minimizer of (D.13), we can write

$$\nabla_{\boldsymbol{\theta}} R(\tilde{\boldsymbol{\theta}}, \boldsymbol{r}) = \frac{1}{n} \sum_{t \neq k} (\nabla \ell(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_t, y_t) - \nabla \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)) - \boldsymbol{H}_{-k}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) + \nabla q(\tilde{\boldsymbol{\theta}}) - \nabla q(\boldsymbol{\theta}_{-k}^*)$$

Now, from (D.14), we obtain

(D.52)
$$\nabla_{\boldsymbol{\theta}} R(\tilde{\boldsymbol{\theta}}, \boldsymbol{r}) = \frac{1}{n} \sum_{t \neq k} \nabla \ell(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_t, y_t) - \nabla \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) - \nabla^2 \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

where in the above we have used the fact that, since q is a quadratic function, we have $\nabla q(\tilde{\boldsymbol{\theta}}) - \nabla q(\boldsymbol{\theta}_{-k}^*) - \nabla^2 q(\boldsymbol{\theta}_{-k}^*)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) = 0$. We will now analyze each of the terms above. We first bound the first term (i.e. the sum involving the derivatives of ℓ). We will use the following simple relations for any choice of r, y:

$$\tilde{\boldsymbol{\theta}}^{\mathsf{T}} \boldsymbol{r} - \boldsymbol{y} = \boldsymbol{\theta}_{-k}^{*} \boldsymbol{r} - \boldsymbol{y} + (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}} \boldsymbol{r},$$

$$\| \boldsymbol{J} \tilde{\boldsymbol{\theta}} \|_{\ell_{2}}^{2} = \| \boldsymbol{J} \boldsymbol{\theta}_{-k}^{*} \|_{\ell_{2}}^{2} + 2 \boldsymbol{\theta}_{-k}^{*} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) + (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}),$$

$$(\tilde{\boldsymbol{\theta}}^{\mathsf{T}} \boldsymbol{r} - \boldsymbol{y})^{2} = (\boldsymbol{\theta}_{-k}^{*} \boldsymbol{r} - \boldsymbol{y})^{2} + 2(\boldsymbol{\theta}_{-k}^{*} \boldsymbol{r} - \boldsymbol{y}) \boldsymbol{r}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) + (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}} \boldsymbol{r} \boldsymbol{r}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}),$$

and

$$h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}, y) - h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}, y) = (\tilde{\boldsymbol{\theta}}^\mathsf{T} \boldsymbol{r} - y)^2 \| J \tilde{\boldsymbol{\theta}} \|_{\ell_2}^2 - (\boldsymbol{\theta}_{-k}^* \mathsf{T} \boldsymbol{r} - y)^2 \| J \boldsymbol{\theta}_{-k}^* \|_{\ell_2}^2$$

$$= \nabla h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}, y)^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) + e(\boldsymbol{r}, y)$$

$$= 2(\| J \boldsymbol{\theta}_{-k}^* \|_{\ell_2}^2 \boldsymbol{r}^\mathsf{T} + (\boldsymbol{\theta}_{-k}^* \mathsf{T} \boldsymbol{r} - y)^2 \boldsymbol{J}^\mathsf{T} J \boldsymbol{\theta}_{-k}^* \mathsf{T}) (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) + e(\boldsymbol{r}, y),$$

where the error term e(r, y) can be written as

$$e(\mathbf{r}, y) = (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \nabla_{\boldsymbol{\theta}}^2 h(\boldsymbol{\theta}; \mathbf{r}, y) |_{\boldsymbol{\theta} = \boldsymbol{\theta}(\mathbf{r}, y)} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

and $\theta(r,y) = \zeta \theta_{-k}^* + (1-\zeta)\tilde{\theta}$ for some $\zeta \in [0,1]$ which depends on r and y.

We will also use the Taylor expansion:

$$g'(h(\tilde{\boldsymbol{\theta}};\boldsymbol{r},y_t)) = g'(h(\boldsymbol{\theta}_{-k}^*;\boldsymbol{r}_t,y_t)) + g''(h(\boldsymbol{\theta}_{-k}^*;\boldsymbol{r}_t,y_t))(h(\tilde{\boldsymbol{\theta}};\boldsymbol{r}_t,y_t) - h(\boldsymbol{\theta}_{-k}^*;\boldsymbol{r}_t,y_t)) + \frac{1}{2}g'''(v_t)(h(\tilde{\boldsymbol{\theta}};\boldsymbol{r}_t,y_t) - h(\boldsymbol{\theta}_{-k}^*;\boldsymbol{r}_t,y_t))^2,$$

where v_t is a number between $h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_t, y_t)$ and $h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)$.

From (D.20) and (D.21) we will decompose:

$$\nabla \ell(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_t, y_t) - \nabla \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) - \nabla^2 \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) = \text{Term}_1 + \text{Term}_2 - \text{Term}_3,$$

where the terms are given in (D.53), (D.55), and (D.56). We will bound each of these terms in the following. We have

(D.53)
$$\operatorname{Term}_{1} = 2\left(\nabla h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_{t}, y_{t}) - \nabla h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t})\right) g'(h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t}))$$

where

$$\nabla h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_t, y_t) - \nabla h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)$$

$$=2\left(\left(\tilde{\boldsymbol{\theta}}^{\mathsf{T}}\boldsymbol{r}_{t}-y_{t}\right)\left\|\boldsymbol{J}\tilde{\boldsymbol{\theta}}\right\|_{\ell_{2}}^{2}\boldsymbol{r}_{t}+\left(\tilde{\boldsymbol{\theta}}^{\mathsf{T}}\boldsymbol{r}_{t}-y_{t}\right)^{2}\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}\tilde{\boldsymbol{\theta}}-\left(\boldsymbol{\theta}_{-k}^{\mathsf{*}\mathsf{T}}\boldsymbol{r}_{t}-y_{t}\right)\left\|\boldsymbol{J}\boldsymbol{\theta}_{-k}^{\mathsf{*}}\right\|_{\ell_{2}}^{2}\boldsymbol{r}_{t}+\left(\boldsymbol{\theta}_{-k}^{\mathsf{*}\mathsf{T}}\boldsymbol{r}_{t}-y_{t}\right)^{2}\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{-k}^{\mathsf{*}}\right)$$

And thus,

(D.54)

$$\nabla h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_t, y_t) - \nabla h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)$$

$$= 2\left((\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \boldsymbol{r}_t \| \boldsymbol{J} \boldsymbol{\theta}_{-k}^* \|_{\ell_2}^2 + (\tilde{\boldsymbol{\theta}}^{\mathsf{T}} \boldsymbol{r}_t - y_t) \left(2\boldsymbol{\theta}_{-k}^*^{\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) + (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \right) \right) \boldsymbol{r}_t$$

$$+ 2\left((\tilde{\boldsymbol{\theta}}^{\mathsf{T}} \boldsymbol{r}_t - y)^2 \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) + \left(2(\tilde{\boldsymbol{\theta}}^{\mathsf{T}} \boldsymbol{r}_t - y_t) \boldsymbol{r}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) + (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \boldsymbol{r}_t \boldsymbol{r}_t^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \right) \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J} \tilde{\boldsymbol{\theta}} \right)$$

We also have

(D.55)
$$\operatorname{Term}_{2} = 2 \nabla h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_{t}, y_{t}) g'' \left(h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t}) \right) \left(h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_{t}, y_{t}) - h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t}) \right)$$
$$= 2 \nabla h(\tilde{\boldsymbol{\theta}}; \boldsymbol{r}_{t}, y_{t}) g'' \left(h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t}) \right) \left(\nabla h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t})^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) + e(\boldsymbol{r}_{t}, y_{t}) \right),$$

and

$$\operatorname{Term}_{3} = 2\left(\nabla^{2}h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t})g'(h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t})) + \nabla h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t})(\nabla h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t}))^{\mathsf{T}}g''(h(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}_{t}, y_{t}))\right)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})$$

$$-\frac{1}{2}g'''(v_t)\left(h(\tilde{\boldsymbol{\theta}};\boldsymbol{r}_t,y_t)-h(\boldsymbol{\theta}_{-k}^*;\boldsymbol{r}_t,y_t)\right)^2,$$

After some straight-forward steps, we can write

(D.57)

 $Term_1 + Term_2 - Term_3$

$$= 4g' \left(h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)\right) \left(2\boldsymbol{r}_t^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{\theta}_{-k}^{*\mathsf{T}} \boldsymbol{J}^\mathsf{T} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{r}_t \right.$$

$$\left. + 2\boldsymbol{r}_t^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \left((\boldsymbol{\theta}_{-k}^{*\mathsf{T}} \boldsymbol{r}_t - y_t) \boldsymbol{J}^\mathsf{T} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) + \boldsymbol{r}_t^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{J}^\mathsf{T} \boldsymbol{J} \tilde{\boldsymbol{\theta}} \right) \right.$$

$$+ \left(2(\boldsymbol{\theta}_{-k}^{*\mathsf{T}}\boldsymbol{r}_{t} - y_{t})\boldsymbol{r}_{t}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) + (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}}\boldsymbol{r}_{t}\boldsymbol{r}_{t}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}))\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})\right)$$

$$+ (\boldsymbol{r}_{t}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}))^{2}\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}\tilde{\boldsymbol{\theta}} + \|\boldsymbol{J}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})\|_{\ell_{2}}^{2}\boldsymbol{r}_{t}\right)$$

$$+ 2g''(h(\boldsymbol{\theta}_{-k}^{*};\boldsymbol{r}_{t},y_{t}))\left(\left(\nabla h(\tilde{\boldsymbol{\theta}};\boldsymbol{r}_{t},y_{t}) - \nabla h(\boldsymbol{\theta}_{-k}^{*};\boldsymbol{r}_{t},y_{t})\right)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}}\nabla h(\boldsymbol{\theta}_{-k}^{*};\boldsymbol{r}_{t},y_{t}) + \nabla h(\tilde{\boldsymbol{\theta}};\boldsymbol{r}_{t},y_{t})e(\boldsymbol{r}_{t},y_{t})\right)$$

$$-\frac{1}{2}g'''(v_{t})\left(h(\tilde{\boldsymbol{\theta}};\boldsymbol{r}_{t},y_{t}) - h(\boldsymbol{\theta}_{-k}^{*};\boldsymbol{r}_{t},y_{t})\right)^{2}.$$

The relation (D.57) has itself three different terms. We will now simplify and bound each of the terms above. However, we remark that all the three terms will be bounded in a similar way. Let's consider the first term in the right-hand-side of (D.57). The first part of this term is:

$$4g'(h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)) \times 2\boldsymbol{r}_t^\mathsf{T}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) {\boldsymbol{\theta}_{-k}^*}^\mathsf{T} \boldsymbol{J}^\mathsf{T} \boldsymbol{J}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{r}_t,$$

which can be rewritten as

$$8g'(h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) \boldsymbol{\theta}_{-k}^{*\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} \boldsymbol{J}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{r}_t \boldsymbol{r}_t^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*).$$

Now, by using the fact that the first derivatives of the function g is uniformly bounded, and by some straight-forward usages of the Cauchy-Schwartz inequality, we can easily rewrite the above part as

$$\alpha_{1,t} \boldsymbol{p}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{r}_t \boldsymbol{r}_t^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

where the vector $\mathbf{p} = \boldsymbol{\theta}_{-k}^* \mathbf{T} \mathbf{J}^\mathsf{T} \mathbf{J}$ does not depend on t, and $\alpha_{1,t}$ is a constant. We can further write:

$$|\alpha_{1,t}| \le C$$
 and $\|\boldsymbol{p}\|_{\ell_2} \le \|J\|_{\ell_2}^2 \|\boldsymbol{\theta}_{-k}^*\|_{\ell_2} \le \|\boldsymbol{J}\|_{\ell_2}^4 + \|\boldsymbol{\theta}_{-k}^*\|_{\ell_2}^2$

where C is an absolute constant.

Let us now consider the second part of the first term in (D.57), which is:

$$4g'(h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)) \times 2\boldsymbol{r}_t^{\mathsf{T}}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)(\boldsymbol{\theta}_{-k}^{*\mathsf{T}}\boldsymbol{r}_t - y_t)\boldsymbol{J}^{\mathsf{T}}\boldsymbol{J}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*).$$

We can rewrite this part as

$$8g'(h(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t)) \times (\boldsymbol{\theta}_{-k}^* \mathsf{T} \boldsymbol{r}_t - y_t) \boldsymbol{J}^\mathsf{T} \boldsymbol{J} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \boldsymbol{r}_t (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) = \alpha_{2,t} \boldsymbol{A} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \boldsymbol{r}_t (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

where, the matrix A is the same for all t, and

$$|\alpha_{2,t}| \le C|\boldsymbol{\theta}_{-k}^*|^{\mathsf{T}} \boldsymbol{r}_t - y_t| \text{ and } \|\boldsymbol{A}\|_{\ell_2} \le \|\boldsymbol{J}\|_{\ell_2}^2.$$

In a similar way, one can inspect all the parts of the first term in (D.57) and show that they take the form of one of the following:

•
$$\alpha_{1,t} p_1^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) r_t r_t^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

•
$$\alpha_{2,t} \mathbf{A}_1 (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} r_t (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

•
$$\alpha_{3,t} p_2 (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} r_t r_t^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

$$\bullet \alpha_{4,t} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \boldsymbol{A}_2 (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{r}_t,$$

(D.58)
$$\bullet \alpha_{5,t} ((\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \boldsymbol{r}_t)^2 \boldsymbol{p}_3,$$

where p_1, p_2, p_3, A_1, A_2 do not depend on t, and we have (D.59)

$$|\alpha_{j,t}| \leq C(1 + (\boldsymbol{\theta}_{-k}^* \boldsymbol{\tau}_{t} - y)^2) \text{ and } \max\{|\gamma|, \|\boldsymbol{p}_1\|_{\ell_2}, \|\boldsymbol{p}_2\|_{\ell_2}, \|\boldsymbol{p}_3\|_{\ell_2}, \|\boldsymbol{A}_1\|_{\ell_2}, \|\boldsymbol{A}_2\|_{\ell_2}\} \leq C(1 + \|\boldsymbol{J}\|_{\ell_2}^2).$$

We will now consider the second term in the right-hand-side of (D.57) which can be expanded using (D.22) and (D.54). Again, one can inspect all the parts and show that they take one of the forms in the following (in addition to the forms presented in (D.58)):

(D.60)
$$\bullet \alpha_{3,t} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{r}_t \boldsymbol{r}_t^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \boldsymbol{p}_3,$$

$$\bullet \alpha_{4,t} (\boldsymbol{r}_t^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*))^2 \boldsymbol{r}_t,$$

$$\bullet \alpha_{5,t} \boldsymbol{p}_4^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \boldsymbol{r}_t \boldsymbol{r}_t^\mathsf{T} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*),$$

where, for some positive constant C and even integer D we have (D.61)

$$|\alpha_{j,t}| \le C(1 + (\boldsymbol{\theta}_{-k}^* \mathbf{T} r - y)^D + (\tilde{\boldsymbol{\theta}}^\mathsf{T} r - y)^D) \text{ and } \max\{\|\boldsymbol{p}_3\|_{\ell_2}, \|\boldsymbol{p}_4\|_{\ell_2}\} \le C(1 + \|\boldsymbol{J}\|_{\ell_2}^D + \|\boldsymbol{\theta}_{-k}^*\|_{\ell_2}^D).$$

A similar bounding can be done for the third term in (D.57).

We now claim that the sum of each of the terms in (D.58) and (D.60) over t is at most of order O(polylog (n)/n). Let's consider the first term in (D.58). We can write

$$\frac{1}{n} \left\| \sum_{t} \alpha_{1,t} \boldsymbol{p}_{1}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) \boldsymbol{r}_{t} \boldsymbol{r}_{t}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) \right\|_{\ell_{2}} \leq \frac{1}{n} \sup_{t} \left\{ |\alpha_{1,t}| \right\} \left| \boldsymbol{p}_{1}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) \right| \left\| \sum_{t \neq k} \boldsymbol{r}_{t} \boldsymbol{r}_{t}^{\mathsf{T}} \right\|_{\ell_{2}} \left\| \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right\|_{\ell_{2}}$$

$$= \sup_{t} \left\{ |\alpha_{1,t}| \right\} \left| \boldsymbol{p}_{1}^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*}) \right| \left\| \frac{1}{n} \sum_{t \neq k} \boldsymbol{r}_{t} \boldsymbol{r}_{t}^{\mathsf{T}} \right\|_{\ell_{2}} \left\| \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right\|_{\ell_{2}}$$

Now, from Lemma D.5 it should be clear why the above quantity is small. Informally, and neglecting the polylog factors, the lemma asserts that with high probability the terms $\sup_t\{|\alpha_{1,t}|\}$, and $\left\|\frac{1}{n}\sum_{t\neq k} r_t r_t^\mathsf{T}\right\|_{\ell_2}$ are all O(1); but the term $\left|p_1^\mathsf{T}(\tilde{\theta}-\theta_{-k}^*)\right|$, is O(1/n) as p_1 is a fixed vector (independent of t), and $\left\|\tilde{\theta}-\theta_{-k}^*\right\|_{\ell_2}$ is $O(n^{-\frac{1}{2}})$. As a result, the whole expression is $O(n^{-\frac{3}{2}})$. Formally, it is easy to conclude from Lemma D.5 that for a given $k \in [d]$:

(D.62)

$$\mathbb{P}\left(\frac{1}{n}\left\|\sum_{t}\alpha_{1,t}\boldsymbol{p}_{1}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^{*})\boldsymbol{r}\boldsymbol{r}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^{*})\right\|_{\ell_{2}} \geq v\frac{(\log(d))^{b}}{d^{\frac{3}{2}}}\right) \leq c\exp(-v^{c'}/c) + c\exp\left(-(\log(d))^{2}/c\right),$$

for some absolute constants c, c', c'' > 0.

Let's now consider the second term in (D.58). We can write

$$\begin{aligned} \left\| \frac{1}{n} \sum_{t} \alpha_{2,t} \mathbf{A}_{1} \left(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right)^{\mathsf{T}} \boldsymbol{r}_{t} \left(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right) \right\|_{\ell_{2}} &= \left\| \mathbf{A}_{1} \right\|_{\ell_{2}} \left\| \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right\|_{\ell_{2}} \left\| \left(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right)^{\mathsf{T}} \frac{1}{n} \sum_{t \neq k} \alpha_{2,t} \boldsymbol{r}_{t} \right\|_{\ell_{2}} \\ &= \left\| \mathbf{A}_{1} \right\|_{\ell_{2}} \left\| \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right\|_{\ell_{2}} \left\| \left(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^{*} \right)^{\mathsf{T}} \times \frac{1}{n} \sum_{t \neq k} \alpha_{2,t} \boldsymbol{r}_{t} \right\|_{\ell_{2}}. \end{aligned}$$

Now, note from the first part of Lemma D.5 that the norm of the vector $\frac{1}{n} \sum_{t \neq k} \alpha_{2,t} \boldsymbol{r}_t^\mathsf{T}$ is w.h.p. O(1). Also, this vector is independent from \boldsymbol{r} , and hence from the second part of Lemma D.5 we obtain that w.h.p. $\|(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*)^\mathsf{T} \times \frac{1}{n} \sum_{t \neq k} \alpha_{2,t} \boldsymbol{r}_t\|_{\ell_2}$ is $O(n^{-\frac{1}{2}})$. Consequently, by noting that $\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*\|_{\ell_2}$ is w.h.p $O(n^{-\frac{1}{2}})$ we obtain that the whole expression is w.h.p. $O(n^{-\frac{3}{2}})$. The formal expression would be like the probabilistic expression given in (D.62).

Similarly, for the term $\alpha_{4,t} (r_t^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*))^2 r_t$ we can write

(D.63)
$$\left\| \frac{1}{n} \sum_{t \neq k} \alpha_{4,t} \left(r_t^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \right)^2 r_t \right\|_{\ell_2} \leq \sup_{t} \left\{ \left(r_t^{\mathsf{T}} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_{-k}^*) \right)^2 \right\} \left\| \frac{1}{n} \sum_{t \neq k} \alpha_{4,t} r_t \right\|_{\ell_2}.$$

Now, by part (e) of Lemma D.5 we can easily conclude that

$$\mathbb{P}\left(\sup_{t}\left\{\left(\boldsymbol{r}_{t}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^{*})\right)^{2}\right\} \geq \left(\log(d)\right)^{b}\right) \leq c\exp\left(-\left(\log(d)\right)^{2}/c\right),$$

for absolute constants b, c > 0 that are suitably chosen. A similar conclusion can be made for $\sup_{t} \{|\alpha_{4,t}|\}$ from (D.61) and part (g) of Lemma D.5–i.e.

$$\mathbb{P}\left(\sup_{t} \{|\alpha_{4,t}|\} \ge (\log(d))^b\right) \le c \exp\left(-(\log(d))^2/c\right).$$

As a result, by using the above bounds, as well as (D.63), and part (a) of Lemma D.5 we obtain

$$\mathbb{P}\left(\left\|\frac{1}{n}\sum_{t\neq k}\alpha_{4,t}\left(\boldsymbol{r}_{t}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^{*})\right)^{2}\boldsymbol{r}_{t}\right\|_{\ell_{2}}\geq v\frac{(\log(d))^{b}}{d}\right)\leq c\exp\left(-v^{c'}/c\right)+c\exp\left(-(\log(d))^{2}/c\right),$$

In a similar way as the above, we can show that the sum of all the terms in (D.58) and (D.60) over t have similar bounds. As a result, going back to (D.52), we have shown that the sum can be bounded to give the desired result as in the lemma.

Lemma D.5 For some absolute constants b, c, c' > 0 we have:

(a)

$$\mathbb{P}\left(\left\|\frac{1}{n}\sum_{t\neq k}\boldsymbol{r}_{t}\boldsymbol{r}_{t}^{\mathsf{T}}\right\|_{\ell_{2}}\geq v\right)\leq c\exp\left(-v^{c'}/c\right).$$

(b) Given any sequence of numbers $\{\alpha_t\}_{t=1}^n$ such that $|\alpha_t| \le 1$, we have

$$\mathbb{P}\left(\left\|\frac{1}{n}\sum_{t\neq k}\alpha_{t}\boldsymbol{r}_{t}\right\|_{\ell_{2}}\geq v\right)\leq ce^{-v^{c'}/c}.$$

(c) Given a vector \mathbf{u} , with $\|\mathbf{u}\|_{\ell_2} = 1$, we can write

$$\mathbb{P}\left(|\boldsymbol{u}^{\mathsf{T}}\boldsymbol{r}| \ge v\right) \le ce^{-v/c}.$$

(d) Given a vector \mathbf{u} , with $\|\mathbf{u}\|_{\ell_2} = 1$, which is independent from \mathbf{r} , we have

$$\mathbb{P}\left(\left|\boldsymbol{u}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^{*})\right| \geq \frac{v}{n}\right) \leq ce^{-v/c}.$$

(e) For r generated according to either of the distributions in (D.18), we have

$$\mathbb{P}\left(\|\boldsymbol{r}\|_{\ell_2} \ge v\sqrt{n}\right) \le c \exp(-v^2/c).$$

(f) We further have

$$\mathbb{P}\left(\left|\boldsymbol{r}_{t}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}-\boldsymbol{\theta}_{-k}^{*})\right| \geq \frac{v}{\sqrt{n}}\right) \leq ce^{-v/c}.$$

(g) For a given even integer D > 0 we have:

$$\mathbb{P}\left(\max\left\{(\boldsymbol{\theta_{-k}^*}^\mathsf{T}\boldsymbol{r}-\boldsymbol{y})^D,(\tilde{\boldsymbol{\theta}}^\mathsf{T}\boldsymbol{r}-\boldsymbol{y})^D,\|\boldsymbol{\theta_{-k}^*}\|_{\ell_2}^D\right\} \geq v\right) \leq ce^{-v^{c'}/c}.$$

As a corollary, we have

$$\mathbb{P}\left(\max\left\{\left(\boldsymbol{\theta_{-k}^*}^\mathsf{T}\boldsymbol{r}-\boldsymbol{y}\right)^D, (\tilde{\boldsymbol{\theta}}^\mathsf{T}\boldsymbol{r}-\boldsymbol{y})^D, \|\boldsymbol{\theta_{-k}^*}\|_{\ell_2}^D\right\} \geq (\log(d))^b\right) \leq c \exp\left\{-(\log(d))^2/c\right\}.$$

Proof of this lemma is provided in Section D.2.7.

D.2.3. Bounding
$$|\mathbf{r}^{\mathsf{T}}\tilde{\boldsymbol{\theta}}(\mathbf{r})|$$

Lemma D.6 There exist absolute constants c, c' > 0 such that for every $k \in [n]$ and $v \ge 0$ we have

(D.64)
$$\mathbb{P}\left(|\boldsymbol{r}^{\mathsf{T}}\tilde{\boldsymbol{\theta}}(\boldsymbol{r})| \ge v\right) \le c \exp(-v^{c'}/c),$$

and

$$(D.65) \qquad \mathbb{P}\left(|\boldsymbol{r}^{\mathsf{T}}\boldsymbol{\theta}_{k}^{*}(\boldsymbol{r})| \ge v\right) \le cd \exp\left(-v^{c'}/c\right) + c \exp\left(-(\log(d))^{2}/c\right).$$

Proof This lemma can also be proven similarly to [39] (see Section D.3 in [39]). There are, however, some small differences that we will mention the details here.

For the first part, the proof proceeds in two steps. In this first step, we use Lemma D.1 to obtain

(D.66)
$$\mathbb{P}\left(|\boldsymbol{r}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*}| \geq v\right) \leq c \exp(-v/c).$$

We will now bound the term $|r^{\mathsf{T}}(\tilde{\theta}(r) - \theta_{-k}^*(r))|$. We can write:

$$\mathbb{P}\left(|\boldsymbol{r}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^{*})| \geq v\right) = \mathbb{P}\left(\left|\frac{1}{\sqrt{d}}\boldsymbol{r}^{\mathsf{T}} \times \sqrt{d}(\tilde{\boldsymbol{\theta}}(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^{*}(\boldsymbol{r}))\right| \geq v\right) \\
\leq \mathbb{P}\left(\left\|\tilde{\boldsymbol{\theta}}(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^{*}\right\|_{\ell_{2}} \geq \sqrt{\frac{v}{d}}\right) + \mathbb{P}\left(\left\|\boldsymbol{r}\right\|_{\ell_{2}} \geq \sqrt{vd}\right)$$

Now, the first term above can be bounded using Lemma D.2, and the second can be bounded from Lemma D.5, and thus the proof of the lemma follows.

The proof of the second part is similar to the first part (and we use Lemma D.4).

D.2.4. Bounding
$$\|\boldsymbol{\theta}_{-k}^*\|_{\ell}$$

Lemma D.7 Let θ_{-k}^* be the minimizer of R_{-k} defined in (D.8). There exist absolute constants $c_0, c_1, c_\infty > 0$ such that for any $k \in [n]$:

(D.67)
$$\mathbb{P}\left\{ \|\boldsymbol{\theta}_{-k}^{*}\|_{\ell_{\infty}} \ge c_{\infty} \sqrt{\frac{\log(d)}{d}} \right\} \le 5d^{-c_{0}} + 3e^{-c_{1}n}$$

Proof For convenience we remind the definition of θ_{-k}^* given by $\theta_{-k}^* = \arg\min_{\theta} R_{-k}(\theta)$ where

(D.68)

$$R_{-k}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{k-1} \ell(\boldsymbol{\theta}; \boldsymbol{a}_i, y_i) + \frac{1}{n} \sum_{i=k+1}^{n} \ell(\boldsymbol{\theta}; \boldsymbol{b}_i, y_i) + \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \left\| \frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \lambda_s \frac{\sqrt{\log(d)}}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta})^2,$$

and

$$\ell(\boldsymbol{\theta}; \boldsymbol{r}, y) = (\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y)^{2} + \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}^{2} + 2g\left((\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r} - y)^{2} \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_{2}}^{2}\right),$$

with
$$g(x) = \sqrt{x + \gamma}$$
.

We use a similar strategy as in the proof of Proposition 6.2. Specifically we first bound the last coordinate of θ_{-k}^* . Next, by symmetry we conclude that the same bound holds for all of its coordinates and control its ℓ_{∞} norm by union bounding.

With a slight abuse of notation, we consider a (N+1) dimensional version of the above optimization over $[\theta; u]$ and denote the last coordinate of the optimal solution by \hat{u} . We define

$$r_i = \begin{cases} \sigma(\boldsymbol{W}\boldsymbol{x}_i) & \text{if } i \leq k-1, \\ \boldsymbol{0} & \text{if } i = k, \\ \boldsymbol{f}_i & \text{if } k \leq i \leq n. \end{cases}$$

Also define $e = [a_1, \dots, a_{k-1}, b_{k+1}, \dots, b_n]$ with $a_i = \sigma(\boldsymbol{w}_{N+1}^\mathsf{T} \boldsymbol{x}_i)$ and $b_i = \mu_1 \boldsymbol{w}_{N+1}^\mathsf{T} \boldsymbol{x}_i + \mu_2 z_i$ (with $z_i \sim \mathsf{N}(0,1)$ independent of \boldsymbol{x}_i). From (D.68), \hat{u} can be expressed as

$$\hat{u} = \arg\min_{u} \min_{\boldsymbol{\theta}} R_{-k}([\boldsymbol{\theta}; u])$$

where

$$R_{-k}([\boldsymbol{\theta}; u]) = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r}_{i} + u e_{i} - y_{i})^{2} + \|\boldsymbol{J}\boldsymbol{\theta} + u\boldsymbol{h}\|_{\ell_{2}}^{2} + 2g((\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{r}_{i} + u e_{i} - y_{i})^{2} \|\boldsymbol{J}\boldsymbol{\theta} + u\boldsymbol{h}\|_{\ell_{2}}^{2})$$

(D.69)
$$+ \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda u^2 + \lambda_w \left\| \frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} + \boldsymbol{w}_{N+1} u - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \lambda_s \frac{\sqrt{\log(d)}}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta} + u)^2,$$

Here $\boldsymbol{h} \in \mathbb{R}^{N+1}$ is the last column of the (N+1)-dimensional matrix \boldsymbol{J} . Namely,

(D.70)
$$\boldsymbol{h} \coloneqq \begin{bmatrix} (\boldsymbol{W}\boldsymbol{w}_{N+1}) \odot \left(\frac{\pi - \cos^{-1}(\boldsymbol{W}\boldsymbol{w}_{N+1})}{2\pi}\right) \\ \frac{1}{2} \end{bmatrix}.$$

Let f(u) denote the objective function of u above, i.e., $f(u) = \min_{\theta} R_{-k}([\theta; u])$. We also let θ_* be the minimizing θ in this objective if we set u = 0, i.e., $\theta_* = \min_{\theta} R_{-k}([\theta; 0])$.

Following a similar argument in the proof of Proposition 6.2, we can obtain a lower bound on f(u) by considering a second-order Taylor expansion of f(u) around $[\theta_*, 0]$ and using the strong-convexity of the loss function to arrive at the following upper bound on \hat{u} :

(D.71)
$$|\hat{u}| \leq \frac{1}{\lambda + \lambda_w} \left| \nabla_u R_{-k}([\boldsymbol{\theta}; u])|_{[\boldsymbol{\theta}_*; 0]} \right|.$$

Calculating $\nabla_u R_{-k}([\boldsymbol{\theta};u])|_{[\boldsymbol{\theta}_*;0]}$ we have

$$\nabla_{u}R_{-k}([\boldsymbol{\theta};u])|_{[\boldsymbol{\theta}_{\star};0]} = \frac{1}{n} \sum_{i=1}^{n} \left[2(\boldsymbol{\theta}_{\star}^{\mathsf{T}}\boldsymbol{r}_{i} - y_{i})e_{i} + 2\boldsymbol{h}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{\star} + \frac{2}{\sqrt{(\boldsymbol{\theta}_{\star}^{\mathsf{T}}\boldsymbol{r}_{i} - y_{i})^{2} \|\boldsymbol{J}\boldsymbol{\theta}_{\star}\|_{\ell_{2}}^{2} + \gamma}} \left((\boldsymbol{\theta}_{\star}^{\mathsf{T}}\boldsymbol{r}_{i} - y_{i}) \|\boldsymbol{J}\boldsymbol{\theta}_{\star}\|_{\ell_{2}}^{2} e_{i} + (\boldsymbol{\theta}_{\star}^{\mathsf{T}}\boldsymbol{r}_{i} - y_{i})^{2} \boldsymbol{h}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{\star} \right) \right]$$

$$+2\lambda_w \boldsymbol{w}_{N+1}^{\mathsf{T}} \left(\frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta}_* - \boldsymbol{\beta}\right) + 2\lambda_s \frac{\sqrt{\log(d)}}{d} \mathbf{1}^{\mathsf{T}} \boldsymbol{\theta}_*.$$

We treat each of these terms separately.

Define the event

$$\mathcal{E} \coloneqq \left\{ \frac{1}{n} \sum_{i=1}^{n} y_i^2 < C, \ \frac{1}{\sqrt{d}} \| \boldsymbol{X} \| \le C, \ \| \boldsymbol{W} \| \le C \right\}.$$

Using the concentration bounds for the operator norm of Gaussian matrices and also the tail bound for chi-square random variables, we have that $\mathbb{P}(\mathcal{E}) \ge 1 - 3e^{-cn}$ for some constant c > 0.

Note that by optimality of $\boldsymbol{\theta}_*$, we have $R_{-k}([\boldsymbol{\theta}_*;0]) \leq R_{-k}([\boldsymbol{0};0]) = \frac{1}{n} \sum_{i=1}^n y_i^2$ from which we get $\|\boldsymbol{\theta}_*\|_{\ell_2} \leq C$ and $\|\frac{1}{2}\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\theta}_* - \boldsymbol{\beta}\|_{\ell_2} \leq C$, on the event \mathcal{E} . Therefore,

$$(D.73) 2\lambda_s \frac{\sqrt{\log(d)}}{d} |\mathbf{1}^\mathsf{T} \boldsymbol{\theta}_*| \le 2\lambda_s \frac{\sqrt{\log(d)}}{d} \|\boldsymbol{\theta}_*\|_{\ell_2} \|\mathbf{1}\|_{\ell_2} \le C \sqrt{\frac{\log(d)}{d}}.$$

Also note that \boldsymbol{w}_{N+1} is drawn independently from \boldsymbol{W} , $\boldsymbol{\theta}_*$ and $\boldsymbol{\beta}$. Since $\boldsymbol{w}_{N+1} \sim \mathrm{Unif}(\mathbb{S}^{d-1})$, given $\frac{1}{2}\boldsymbol{W}^\mathsf{T}\boldsymbol{\theta}_* - \boldsymbol{\beta}$ the conditional distribution of $\boldsymbol{w}_{N+1}^\mathsf{T}(\frac{1}{2}\boldsymbol{W}^\mathsf{T}\boldsymbol{\theta}_* - \boldsymbol{\beta})$ converges to $\mathsf{N}(0,\frac{1}{d}\|\frac{1}{2}\boldsymbol{W}^\mathsf{T}\boldsymbol{\theta}_* - \boldsymbol{\beta}\|_{\ell_2}^2)$ form which we obtain

(D.74)
$$\left| \boldsymbol{w}_{N+1}^{\mathsf{T}} \left(\frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta}_{*} - \boldsymbol{\beta} \right) \right| \leq \sqrt{2c' \frac{\log(d)}{d}} \left\| \frac{1}{2} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{\theta}_{*} - \boldsymbol{\beta} \right\|_{\ell_{2}} < C \sqrt{2c' \frac{\log(d)}{d}} ,$$

with probability at least $1 - d^{-c'}$.

We next focus on the terms in the right-hand side of (D.72), which involve e_i . This part can be written as $\frac{1}{\sqrt{n}} \sum_{i=1}^{n} m_i e_i$ with

$$m_i = \frac{1}{\sqrt{n}} \left[2(\boldsymbol{\theta}_*^\mathsf{T} \boldsymbol{r}_i - y_i) + \frac{2}{\sqrt{(\boldsymbol{\theta}_*^\mathsf{T} \boldsymbol{r}_i - y_i)^2 \|\boldsymbol{J} \boldsymbol{\theta}_*\|_{\ell_2}^2 + \gamma}} (\boldsymbol{\theta}_*^\mathsf{T} \boldsymbol{r}_i - y_i) \|\boldsymbol{J} \boldsymbol{\theta}_*\|_{\ell_2}^2 \right].$$

Note that

$$|m_i| \le \frac{2}{\sqrt{n}} \left(|\boldsymbol{\theta}_*^\mathsf{T} \boldsymbol{r}_i - y_i| + ||\boldsymbol{J} \boldsymbol{\theta}_*||_{\ell_2} \right).$$

By optimality of θ_* , on the event \mathcal{E} we have

$$\|\boldsymbol{m}\|_{\ell_2}^2 \le \mathbb{R}_{-k}([\boldsymbol{\theta}_*;0]) \le \mathbb{R}_{-k}([\boldsymbol{0};0]) = \frac{1}{n} \sum_{i=1}^n y_i^2 < C.$$

Observe that w_{N+1} is independent from $\{m_i\}_{i\in[n]}$ (recall that θ_* does not depend on w_{N+1} by its definition.) Following the same strategy in the proof of Proposition 6.2, we only consider the randomness in w_{N+1} and condition on everything else. Write $\frac{1}{\sqrt{n}}\sum_{i=1}^n m_i e_i$ as a function of w_{N+1} as follows:

(D.75)
$$V(\boldsymbol{w}_{N+1}) := \frac{1}{\sqrt{n}} \sum_{i=1}^{k-1} m_i \sigma(\boldsymbol{w}_{N+1}^\mathsf{T} \boldsymbol{x}_i) + \frac{1}{\sqrt{n}} \sum_{i=k+1}^n m_i (\mu_1 \boldsymbol{w}_{N+1}^\mathsf{T} \boldsymbol{x}_i + \mu_2 z_i).$$

Observe that the (conditional) expectation $\mathbb{E}[V(\boldsymbol{w}_{N+1})|\boldsymbol{W},\boldsymbol{X}]=0$. In addition, $V(\cdot)$ is a Lipschitz function with Lipschitz factor at most $\frac{C}{\sqrt{d}}\|\boldsymbol{X}\boldsymbol{m}\|_{\ell_2}$. Therefore, using the concentration bound for Lipschitz function on unit sphere (see e.g. [92, Theorem 5.1.4]), we obtain

$$\mathbb{P}\left(|V(\boldsymbol{w}_{N+1})| \ge t\right) \le \mathbb{P}\left(|V(\boldsymbol{w}_{N+1})| \ge t; \mathcal{E}\right) + \mathbb{P}(\mathcal{E}^c)$$

$$\le 2e^{-c'dt^2} + 3e^{-cn}.$$

Choosing $t = C\sqrt{\frac{\log(d)}{d}}$, we get

(D.76)
$$\mathbb{P}\left(|V(\boldsymbol{w}_{N+1})| \ge C\sqrt{\frac{\log(d)}{d}}\right) \le 2d^{-c'C^2} + 3e^{-cn}.$$

The remaining terms in (D.72) can be rearranged and written as $Ah^{\mathsf{T}}J\theta_{\star}$, with

$$A \coloneqq 2 \left(1 + \frac{1}{n} \sum_{i=1}^{n} \frac{(\boldsymbol{\theta}_* \boldsymbol{r}_i - y_i)^2}{\sqrt{(\boldsymbol{\theta}_*^\mathsf{T} \boldsymbol{r}_i - y_i)^2 \|\boldsymbol{J} \boldsymbol{\theta}_*\|_{\ell_2}^2 + \gamma}} \right).$$

We next bound A > 0:

$$|A| < 2\left(1 + \frac{1}{n} \sum_{i=1}^{n} \frac{|\boldsymbol{\theta}_{*} \boldsymbol{r}_{i} - y_{i}|}{\|\boldsymbol{J} \boldsymbol{\theta}_{*}\|_{\ell_{2}}}\right)$$

$$2\left(1 + \frac{1}{\|\boldsymbol{J} \boldsymbol{\theta}_{*}\|_{\ell_{2}}} \left(\frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{\theta}_{*}^{\mathsf{T}} \boldsymbol{r}_{i} - y_{i})^{2}\right)^{1/2}\right)$$

$$< 2\left(1 + \frac{1}{\|\boldsymbol{J} \boldsymbol{\theta}_{*}\|_{\ell_{2}}} \left(R_{-k}([\boldsymbol{\theta}_{*}; 0])\right)^{1/2}\right)$$

$$\leq 2\left(1 + \frac{1}{\|\boldsymbol{J} \boldsymbol{\theta}_{*}\|_{\ell_{2}}} \left(R_{-k}([\boldsymbol{0}; 0])\right)^{1/2}\right)$$

$$= 2\left(1 + \frac{1}{\|\boldsymbol{J} \boldsymbol{\theta}_{*}\|_{\ell_{2}}} \left(\frac{1}{n} \sum_{i=1}^{n} y_{i}^{2}\right)^{1/2}\right).$$

Using Lemma F.3, on the event \mathcal{E} , the right-hand side of the above equation is of order one (A < C), for some constant C > 0).

We next bound $h^{\mathsf{T}}J\theta_*$. Using the relation $\frac{1}{2\pi}(\pi - \cos^{-1}(\rho)) = \frac{1}{4} + \frac{\rho}{2\pi} + O(\rho^3)$, we define \tilde{h} as follows

$$ilde{m{h}} \coloneqq \begin{bmatrix} rac{1}{4} m{W} m{w}_{n+1}^{\mathsf{T}} \\ rac{1}{2} \end{bmatrix}.$$

Recalling h given by (D.70), we have $\|h - \tilde{h}\|_{\ell_2} = O(1/\sqrt{d})$. On the event \mathcal{E} , we have $\|\theta_*\|_{\ell_2} = O(1)$. Also by invoking Lemma F.3, on event \mathcal{E} , we have $\|J\| = O(1)$. Hence,

(D.77)
$$|\boldsymbol{h}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{\star} - \tilde{\boldsymbol{h}}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{\star}| \leq O(1/\sqrt{d}).$$

We henceforth focus on bounding $\tilde{\boldsymbol{h}}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_*$. Recall that \boldsymbol{w}_{N+1} is independent of \boldsymbol{W} and $\boldsymbol{\theta}_*$. Viewing $\tilde{\boldsymbol{h}}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_*$ as a function of \boldsymbol{w}_{N+1} , it has zero expectation (w.r.t \boldsymbol{w}_{N+1} conditioned on $\boldsymbol{\theta}_*$ and \boldsymbol{W}). In addition, it is a Lipschitz function with Lipschitz factor at most $\frac{1}{4}\|\boldsymbol{W}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_*\|_{\ell_2}$, which is O(1) on the event \mathcal{E} . Next, by employing the concentration bound for Lipschitz functions on unit sphere (see e.g. [92, Theorem 5.1.4]), we obtain

$$\mathbb{P}\left(|\tilde{\boldsymbol{h}}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{*}| \geq t\right) \leq \mathbb{P}\left(|\tilde{\boldsymbol{h}}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{*}| \geq t; \mathcal{E}\right) + \mathbb{P}(\mathcal{E}^{c})$$

$$\leq 2e^{-c'dt^{2}} + 3e^{-cn}.$$

Choosing $t = C\sqrt{\frac{\log(d)}{d}}$, and invoking (D.77) we get

(D.78)
$$\mathbb{P}\left(A\left|\boldsymbol{h}^{\mathsf{T}}\boldsymbol{J}\boldsymbol{\theta}_{*}\right| \geq C\sqrt{\frac{\log(d)}{d}}\right) \leq 2d^{-c'C^{2}} + 3e^{-cn}.$$

Combining the bounds (D.73) to (D.78) into (D.71) we get

$$|\hat{u}| \le C' \frac{\log(d)}{\sqrt{d}},$$

with probability at least $4d^{-c^{\prime}C^2} + 3e^{-cn} + d^{-c^{\prime}}$.

The result follows by union bounding over the N coordinates of $\widehat{\theta}$, along with the assumption that N, n, d grow at the same order. (Note that the event \mathcal{E} is common across all these bounds and so we count its complement probability once.)

D.2.5. Bounding the Difference of Φ_k , Φ_{k-1} with $\Psi_k(r)$

Lemma D.8 We have

(D.79)
$$\max \left\{ \mathbb{E}[(\Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k)) - \Phi_{-k})^2], (\Psi_k(\boldsymbol{f}_k) - \Phi_{-k})^2 \right\} \leq \frac{\text{polylog}(d)}{d^2}.$$

and

(D.80)
$$\max \left\{ \mathbb{E}\left[\left(\Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k)) - \Phi_{k-1} \right)^2 \right], \left(\Psi_k(\boldsymbol{f}_k) - \Phi_k \right)^2 \right\} \le \frac{\text{polylog}(d)}{d^3}$$

Proof To prove (D.79), we note from (D.13) that

$$\Psi_{k}(\boldsymbol{r}) - \Phi_{-k} = \min_{\boldsymbol{\theta}} S_{k}(\boldsymbol{\theta}, \boldsymbol{r}) \leq S_{k}(\boldsymbol{\theta}_{-k}^{*}, \boldsymbol{r})$$

$$= \frac{1}{n} \ell(\boldsymbol{\theta}_{-k}^{*}; \boldsymbol{r}, y_{k})$$

$$\leq \frac{C}{n} \left(1 + |y_{k}| + ||\boldsymbol{J}\boldsymbol{\theta}_{-k}^{*}||_{\ell_{2}} + \boldsymbol{r}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*} \right)^{2},$$

where C is an absolute constant. Consequently, by using Lemma D.1 and the fact that $\|J\|$ is bounded, we obtain (D.79).

To prove (D.80), we adapt the proof of Lemma 1 in [39] to our setting. In the following we use

$$q(\boldsymbol{\theta}) = \lambda \|\boldsymbol{\theta}\|_{\ell_2}^2 + \lambda_w \left\| \frac{1}{2} \boldsymbol{W}^\mathsf{T} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_2}^2 + \lambda_s \frac{\log(d)}{d} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta})^2$$

We start with writing the Taylor expansion of $R_k(\theta, r)$, defined in (D.9), around the point θ_{-k}^* . Note that θ_{-k}^* is the minimizer of $R_{-k}(\theta)$, and hence

$$R_{k}(\boldsymbol{\theta}, \boldsymbol{r}) = R_{-k}(\boldsymbol{\theta}_{-k}^{*}) + \frac{1}{n}\ell(\boldsymbol{\theta}; \boldsymbol{r}, y_{k}) + \frac{1}{2n}\sum_{t \neq k}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}}\nabla^{2}\ell(\boldsymbol{\theta}'; \boldsymbol{r}_{t}, y_{t})(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*}) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})^{\mathsf{T}}\nabla^{2}q(\boldsymbol{\theta}')(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^{*})$$

where θ' can be written as

(D.81)
$$\boldsymbol{\theta}' = \omega \boldsymbol{\theta}_{-k}^* + (1 - \omega) \boldsymbol{\theta},$$

for some $\omega \in [0,1]$. As a result, we can write using the definition (D.13)

$$R_k(\boldsymbol{\theta}, \boldsymbol{r}) - S_k(\boldsymbol{\theta}, \boldsymbol{r})$$

(D.82)
$$= \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \left[\frac{1}{n} \sum_{t \neq k} \nabla^2 \ell(\boldsymbol{\theta}'; \boldsymbol{r}_t, y_t) - \nabla^2 \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*),$$

where we have noted that, since q is a quadratic function, we have $\nabla^2 q(\boldsymbol{\theta}_{-k}^*) = \nabla^2 q(\boldsymbol{\theta}')$.

Let us now consider the sum involving the terms of the form

(D.83)
$$(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \left(\nabla^2 \ell(\boldsymbol{\theta}'; \boldsymbol{r}_t, y_t) - \nabla^2 \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) \right) (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*).$$

We can now use the expansion in (D.21) to bound the above term. A straight-forward calculation, similar to what was done in the proof of Lemma D.4, shows that the above term involves several terms, among which the dominant term has the following form:

$$\alpha_t(\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \left(\boldsymbol{r}_t^{\mathsf{T}} (\boldsymbol{r}_t^{\mathsf{T}} (\boldsymbol{\theta}' - \boldsymbol{\theta}_{-k}^*)) \right) (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) = \alpha_t \left(\boldsymbol{r}_t^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) \right)^2 \left(\boldsymbol{r}_t^{\mathsf{T}} (\boldsymbol{\theta}' - \boldsymbol{\theta}_{-k}^*) \right) = \alpha_t \omega \left(\boldsymbol{r}_t^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) \right)^3,$$

where α_t satisfies (noting that the derivatives of g are uniformly bounded):

(D.85)
$$\mathbb{P}(|\alpha_t| \ge v) \le c \exp(-v^{c'}/c),$$

for absolute constants c, c' > 0. A straight-forward calculation (similar to what is done in the proof of Lemma D.4) shows that all the other terms in the expansion of (D.83) are in absolute value less than the term given in (D.84). As a result, one can write

$$\left| \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*)^{\mathsf{T}} \left[\frac{1}{n} \sum_{t \neq k} \nabla^2 \ell(\boldsymbol{\theta}'; \boldsymbol{r}_t, y_t) - \nabla^2 \ell(\boldsymbol{\theta}_{-k}^*; \boldsymbol{r}_t, y_t) \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) \right| \leq \frac{1}{n} \sum_{t \neq k} |\alpha_t| \left| \boldsymbol{r}_t^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) \right|^3,$$

Using the above bound, we can now bound (D.82) as

(D.86)
$$|R_k(\boldsymbol{\theta}, \boldsymbol{r}) - S_k(\boldsymbol{\theta}, \boldsymbol{r})| \le \frac{1}{n} \sum_{t+k} |\alpha_t| \left| \boldsymbol{r}_t^{\mathsf{T}} (\boldsymbol{\theta} - \boldsymbol{\theta}_{-k}^*) \right|^3.$$

The rest of the proof follows almost line-by-line according to the proof of Lemma 1 in [39]. Let $\mathcal{B} = \{\theta_k^*(r)\} \cup \{\tilde{\theta}_k(r)\}$. By using the definitions (D.12) and (D.16), we have

$$|\Phi_k(r) - \Psi_k(r)| = \left| \min_{\theta \in \mathcal{B}} R_k(\theta, r) - \min_{\theta \in \mathcal{B}} S_k(\theta, r) \right|$$

$$\leq \max_{\theta \in \mathcal{B}} |R_k(\theta, r) - S_k(\theta, r)|$$

We thus obtain using (D.86) that

$$|\Phi_k(\boldsymbol{r}) - \Psi_k(\boldsymbol{r})| \le C \frac{1}{n} \sum_{t \neq k} |\alpha_t| \left(\left| \boldsymbol{r}_t^\mathsf{T} (\boldsymbol{\theta}_k^*(\boldsymbol{r}) - \tilde{\boldsymbol{\theta}}_k(\boldsymbol{r})) \right|^3 + \left| \boldsymbol{r}_t^\mathsf{T} (\tilde{\boldsymbol{\theta}}_k(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^*) \right|^3 \right)$$

Let us now bound each of the terms above. We have

$$\frac{1}{n} \sum_{t \neq k} |\alpha_t| \left| \boldsymbol{r}_t^\mathsf{T} (\boldsymbol{\theta}_k^*(\boldsymbol{r}) - \tilde{\boldsymbol{\theta}}_k(\boldsymbol{r})) \right|^3 \leq \left\| \boldsymbol{\theta}_k^*(\boldsymbol{r}) - \tilde{\boldsymbol{\theta}}_k(\boldsymbol{r}) \right\|_{\ell_2}^3 \frac{1}{n} \sum_{t \neq k} |\alpha_t| \left\| \boldsymbol{r}_t \right\|_{\ell_2}^3.$$

Now, from Lemma D.3, up to negligible $O(\frac{\text{polylog}(n)}{n})$ terms, we have

$$\boldsymbol{r}_t^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}_k(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^*) = \frac{1}{n}\beta_1 \boldsymbol{r}_t^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \boldsymbol{r} + \frac{1}{n}\beta_2 \boldsymbol{r}_t^{\mathsf{T}} H_{-k}^{-1} \boldsymbol{p}.$$

Now, using the above relations, and the inequality $(\sum_{i=1}^{n} |a_i|)^2 \le n \sum_{i=1}^{n} a_i^2$, as well as the Holder's inequality, we can write

$$\begin{aligned} \left| \Phi_{k}(\boldsymbol{r}) - \Psi_{k}(\boldsymbol{r}) \right|^{2} &\leq C'' \left\| \boldsymbol{\theta}_{k}^{*}(\boldsymbol{r}) - \tilde{\boldsymbol{\theta}}_{k}(\boldsymbol{r}) \right\|_{\ell_{2}}^{6} \left(\frac{1}{n} \sum_{t \neq k} |\alpha_{t}|^{2} \left\| \boldsymbol{r}_{t} \right\|_{\ell_{2}}^{6} \right) \\ &+ \frac{C''}{n} \sum_{t \neq k} |\alpha_{t}|^{2} \left(\left| \beta_{1} \frac{\boldsymbol{r}_{t}^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \boldsymbol{r}}{n} \right|^{6} + \left| \beta_{2} \frac{\boldsymbol{r}^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \boldsymbol{p}}{n} \right|^{6} \right) \\ &+ O\left(\frac{\text{polylog}(n)}{n^{4}} \right), \end{aligned}$$

where C'' > 0 is an absolute constant. Now, from Lemma D.4, parts (c) and (e) of Lemma D.5, and (D.34), and (D.85), we obtain for any integer D > 0 that

$$\mathbb{E}\left[\left\|\boldsymbol{\theta}_{k}^{*}(\boldsymbol{r}) - \tilde{\boldsymbol{\theta}}_{k}(\boldsymbol{r})\right\|_{\ell_{2}}^{2D}\right] \leq C_{D} \frac{\operatorname{polylog}(n)}{n^{2D}},$$

$$\mathbb{P}\left(\left\|\boldsymbol{r}_{t}\right\|_{\ell_{2}}^{2D} \geq n^{D}(\log(n))^{b}\right) \leq c \exp\left(-(\log(d))^{2}/c\right),$$

$$\mathbb{P}\left(\left|\alpha_{t}\right|^{2D} \geq (\log(n))^{b}\right) \leq c \exp\left(-(\log(d))^{2}/c\right),$$

$$\mathbb{P}\left(\max\{\left|\beta_{1}\right|^{2D}, \left|\beta_{2}\right|^{2D}\} \geq (\log(n))^{b}\right) \leq c \exp\left(-(\log(d))^{2}/c\right),$$

for suitably chosen absolute constants $b, c, C_D > 0$. Finally, since the matrix \mathbf{H}_{-k}^{-1} has bounded norm, and \mathbf{r}_t and \mathbf{r} are independent sub-gaussian random vectors, we obtain

$$\mathbb{E}\left[\left|\frac{\boldsymbol{r}_{t}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}}{n}\right|^{2D}\right] \leq C_{D}\frac{\operatorname{polylog}(n)}{n^{D}}$$
and
$$\mathbb{E}\left[\left|\frac{\boldsymbol{r}_{t}\boldsymbol{H}_{-k}^{-1}\boldsymbol{p}}{n}\right|^{2D}\right] \leq \mathbb{E}\left[\left|\frac{\|\boldsymbol{r}_{t}\|_{\ell_{2}}\|\boldsymbol{H}_{-k}^{-1}\|_{\ell_{2}}\|\boldsymbol{p}\|_{\ell_{2}}}{n}\right|^{2D}\right] \leq C_{D}\frac{\operatorname{polylog}(n)}{n^{D}}.$$

By using the above relations, the following result now follows in a straight-forward manner using the Holder's inequality:

$$\mathbb{E}\Big[\left|\Phi_k(\boldsymbol{r}) - \Psi_k(\boldsymbol{r})\right|^2\Big] \leq \frac{\operatorname{polylog}(n)}{n^3}.$$

And the result follows since d and n grow in proportion to each other.

D.2.6. Putting things together To prove Theorem 6.9, we consider any test function $\phi: \mathbb{R} \to \mathbb{R}$ which is uniformly bounded in terms of its value as well as its first and second derivatives. We will show that

(D.87)
$$|\mathbb{E}[\varphi(\Phi_A)] - \mathbb{E}[\varphi(\Phi_B)]| = \frac{\text{polylog}(d)}{d^{\frac{1}{2}}} + o_d(1).$$

Using this result, one immediately obtains the theorem (see Sections 2.3 and 2.4 of [39]). As a result, in the rest of this section we focus on proving the above relation for any test function φ . In order to prove this result, we use the so-called Lindeberg's method: We consider the quantities Φ_k defined in (D.7), and show for any $k \in [n]$ that

(D.88)
$$|\mathbb{E}[\varphi(\Phi_k)] - \mathbb{E}[\varphi(\Phi_{k-1})]| = \frac{\text{polylog}(d)}{d^{\frac{3}{2}}} + \frac{o_d(1)}{d}.$$

The above bound immediately results in (D.87) via a telescopic sum over k. It thus remains to prove (D.88).

Using the Taylor expansion, we can write

$$\varphi(\Phi_k) = \varphi(\Phi_{-k}) + \varphi'(\Phi_{-k})(\Phi_k - \Phi_{-k}) + \frac{1}{2}\varphi''(\alpha)(\Phi_k - \Phi_{-k})^2,$$

where α is a number between Φ_{-k} and Φ_k . Using the above expansion, and a similar expansion for Φ_{k-1} , we obtain (D.89)

$$|\mathbb{E}[\varphi(\Phi_k)] - \mathbb{E}[\varphi(\Phi_{k-1})]| \le ||\varphi'||_{\infty} |\mathbb{E}[\Phi_k - \Phi_{k-1}]| + \frac{1}{2} ||\varphi''||_{\infty} ((\Phi_k - \Phi_{-k})^2 + (\Phi_{k-1} - \Phi_{-k})^2),$$

where $||\varphi'||_{\infty}$ and $||\varphi''||_{\infty}$ are the maximum (absolute) values of the first and second derivative of φ .

By using Lemma D.8 we obtain that

$$\begin{aligned} |\mathbb{E}[\Phi_k - \Phi_{k-1}]| &\leq |\mathbb{E}[\Psi_k(\boldsymbol{f}_k) - \Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k))]| + \mathbb{E}[|\Phi_k - \Psi_k(\boldsymbol{f}_k)|] + \mathbb{E}[|\Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k)) - \Phi_{k-1}|] \\ &\leq |\mathbb{E}[\Psi_k(\boldsymbol{f}_k) - \Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k))]| + \frac{\text{polylog}(d)}{d^{\frac{3}{2}}}, \end{aligned}$$

where the last step follows simply from (D.80). Also, from (D.79) and (D.80), we can conclude that

$$\mathbb{E}[(\Phi_k - \Phi_{-k})^2] \le 2\mathbb{E}[(\Phi_k - \Psi_k(\boldsymbol{f}_k))^2] + 2\mathbb{E}[(\Phi_{-k} - \Psi_k(\boldsymbol{f}_k))^2] \le \frac{\text{polylog}(d)}{d^2},$$

and similarly

$$\mathbb{E}[(\Phi_{k-1} - \Phi_{-k})^2] \le \frac{\text{polylog}(d)}{d^2}.$$

Finally, the only term that is left to be analyzed is $|\mathbb{E}[\Psi_k(f_k) - \Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k))]|$, for which we use Lemma D.3, D.7, as well as a CLT-type result from [30]. We note the following three facts:

(i) The quantity $r^{\mathsf{T}}\theta_{-k}^*$ converges in distribution to a gaussian with the same mean and variance when r is generated according to the distributions in (D.18). This is due to the CLT-type theorem given in [30, Theorem 2]. More precisely, we have shown in Lemma D.7 that with probability $1 - cd^{-c}$ we have: $\|\theta_{-k}^*\|_{\ell_{\infty}}$ is at most $C\sqrt{(\log(d))/d}$, where c, C are absolute constants. Also, according to part (e) of Lemma D.1, we have $|\mathbf{1}^{\mathsf{T}}\theta_{-k}^*| \leq C'\sqrt{d/(\log(d))}$, where C' is an absolute constant.

Now, let us define $\theta' = \theta_{-k}^* / \sqrt{\log(d)}$. Note that, with high probability (as specified above), we have $\|\theta'\|_{\ell_{\infty}} \leq C / \sqrt{d}$ and $|\mathbf{1}^{\mathsf{T}}\theta'| \leq C' \sqrt{d} / (\log(d))$. According to [30, Theorem 2], for a and b generated according (D.18), and fixing θ' , the random variables $a^{\mathsf{T}}\theta'$ and $b^{\mathsf{T}}\theta'$ have the same mean and variance, and we have

$$d_{\mathrm{MS}}\left(\boldsymbol{a}^{\mathsf{T}}\boldsymbol{\theta}',\boldsymbol{b}^{\mathsf{T}}\boldsymbol{\theta}'\right) \leq C'' \left\|\boldsymbol{\theta}'\right\|_{\ell_{\infty}} \left(\frac{\left|\mathbf{1}^{\mathsf{T}}\boldsymbol{\theta}'\right|}{\sqrt{d}} + \frac{1}{\sqrt{d}}\right),$$

where C'' > 0 is an absolute constant, and $d_{\rm MS}$ is the so-called maximum-sliced distance, and $d_{\rm MS}\left(\boldsymbol{a}^{\mathsf{T}}\boldsymbol{\theta}',\boldsymbol{b}^{\mathsf{T}}\boldsymbol{\theta}'\right)$ defines the distance between the distributions of $\boldsymbol{a}^{\mathsf{T}}\boldsymbol{\theta}'$ and $\boldsymbol{b}^{\mathsf{T}}\boldsymbol{\theta}'$. As a result, since $\boldsymbol{\theta}_{-k}^* = \boldsymbol{\theta}' \times \sqrt{\log(d)}$, we obtain that

$$d_{\mathrm{MS}}\left(\boldsymbol{a}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*},\boldsymbol{b}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*}\right) \leq C''\sqrt{\log(d)} \left\|\boldsymbol{\theta}'\right\|_{\ell_{\infty}} \left(\frac{\left|\mathbf{1}^{\mathsf{T}}\boldsymbol{\theta}'\right|}{\sqrt{d}} + \frac{1}{\sqrt{d}}\right) = O\left(\frac{1}{\sqrt{\log(d)}}\right).$$

(ii) Consider the result of Lemma D.3. From Lemma D.1, the norm of the vector $\boldsymbol{\theta}_{-k}^*$ is bounded by an absolute constant with probability at least $1-\exp(-cn)$. Hence, since the matrix \boldsymbol{J} is also of bounded operator norm, then the norm of the vector \boldsymbol{p} given in Lemma D.3 is bounded by an absolute constant. Also, the quantity $\|\boldsymbol{J}\boldsymbol{\theta}_{-k}^*\|_{\ell_2}$ is bounded by an absolute constant. Given fixed matrix \boldsymbol{H}^{-1} with bounded norm, and a fixed vector \boldsymbol{p} with bounded norm, the quantity $\frac{1}{n}\boldsymbol{r}^{\mathsf{T}}\boldsymbol{H}^{-1}\boldsymbol{p}$ is, with probability at least $1-c\exp(-(\log(d))^2/c)$, of order $O(\operatorname{polylog}(d)/d)$ according to Lemma D.9. Hence, in the formula (D.31), the overall contribution of the terms which include $\frac{1}{n}\boldsymbol{r}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{p}$ or $\frac{1}{n}\boldsymbol{p}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{p}$ is of order $O(\operatorname{polylog}(d)/d^2)$. Therefore, neglecting these terms adds an additional error of at most $O(\operatorname{polylog}(d)/d^2)$ in computing $\mathbb{E}[\Psi(\boldsymbol{r})] - \Phi_{-k}$. Consequently, from the result of Lemma D.3 we can write

$$\mathbb{E}[\Psi_k(\boldsymbol{r})] = \Phi_{-k} + \frac{1}{n} \mathbb{E} \left[\min_{\tau_1} \left\{ \frac{\boldsymbol{r}^\mathsf{T} \boldsymbol{H}_{-k}^{-1} \boldsymbol{r}}{2n} \left(\frac{\partial \tilde{\ell}(\tau_1, 0)}{\partial \tau_1} \right)^2 + \tilde{\ell}(\tau_1, 0) \right\} \right] + O\left(\frac{\mathsf{polylog}(d)}{d^{\frac{3}{2}}} \right),$$

where $\tilde{\ell}(\tau_1, \tau_2)$ is given in (D.32).

(iii) Given a matrix H, the value $\frac{1}{n}r^{\mathsf{T}}H_{-k}^{-1}r$ concentrates on the same quantity if r is generated from either of the distributions in (D.18). More precisely, from [39, Lemma 13] (or [56, Lemma 1]) we obtain

$$\mathbb{P}\left(\left|\frac{1}{n}\boldsymbol{r}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r} - \mathbb{E}\left[\frac{1}{n}\boldsymbol{r}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right]\right| \ge c\frac{\log(d)}{\sqrt{d}}\right) \le 1 - c\exp\left(-(\log(d)^2)\right),$$

for r being generated according to either of the distributions in (D.18), and c > 0 being an absolute constant. Also, we have that

$$\left| \mathbb{E} \left[\frac{1}{n} \sigma(\boldsymbol{W} \boldsymbol{x})^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \sigma(\boldsymbol{W} \boldsymbol{x}) \right] - \mathbb{E} \left[\frac{1}{n} \boldsymbol{f}^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \boldsymbol{f} \right] \right| = \frac{1}{n} \operatorname{Trace} \left(\boldsymbol{H}_{-k}^{-1} (\Sigma_s - \Sigma_f) \right),$$

where $\Sigma_2 = \mathbb{E}[\sigma(\boldsymbol{W}\boldsymbol{x})\sigma(\boldsymbol{W}\boldsymbol{x})^\mathsf{T}]$ and $\mathbb{E}[\boldsymbol{f}\boldsymbol{f}^\mathsf{T}]$, and the last inequality follows from Lemma D.10. As a result, we obtain

$$\mathbb{P}\left(\left|\frac{1}{n}\sigma(\boldsymbol{W}\boldsymbol{x}_{k})^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\sigma(\boldsymbol{W}\boldsymbol{x}_{k}) - \frac{1}{n}\mathbb{E}\left[\boldsymbol{f}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{f}\right]\right| \ge c\frac{(\log(d))^{3/2}}{\sqrt{d}}\right) \le 1 - c\exp(-(\log(d)^{2})),$$

$$\mathbb{P}\left(\left|\frac{1}{n}\boldsymbol{f}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{f} - \frac{1}{n}\mathbb{E}\left[\boldsymbol{f}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{f}\right]\right| \ge c\frac{\log(d)}{\sqrt{d}}\right) \le 1 - c\exp(-(\log(d)^{2})),$$

Let us now put all the above facts together to bound $|\mathbb{E}[\Psi_k(f_k) - \Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k))]|$. Consider the function $\tilde{\ell}(\tau_1, \tau_2)$ given in (D.32). This function depends on \boldsymbol{r} only through $\boldsymbol{\theta}_{-k}^*\boldsymbol{r}$. Using fact (i) above, we know that $\boldsymbol{\theta}_{-k}^*\boldsymbol{r}$ will asymptotically have the same (gaussian) distribution for both $\boldsymbol{r} \sim f_k$ and $\boldsymbol{r} \sim \sigma(\boldsymbol{W}\boldsymbol{x}_k)$. Also, it is easy to conclude using part (c) of Lemma D.1 that all of the moments of random variable $\boldsymbol{\theta}_{-k}^*\boldsymbol{r}$ are bounded (i.e. $\mathbb{E}[|\boldsymbol{\theta}_{-k}^*\boldsymbol{r}|^D] \leq C_D$ for an absolute constant $C_D > 0$). Further, in fact (ii) we have argued that $\|J\boldsymbol{\theta}_{-k}^*\|_{\ell_2}$ is bounded with probability $1 - e\exp(-cn)$. Also, from fact (iii) above, we know that the term $\frac{1}{n}\boldsymbol{r}^T\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}$ concentrates sharply on the same value for \boldsymbol{r} being either \boldsymbol{f}_k or $\sigma(\boldsymbol{W}\boldsymbol{x}_k)$. Putting all these together, and using [7, Corollary of Theorem 25.12], we obtain that

$$\mathbb{E}_{\boldsymbol{r}=\boldsymbol{f}_{k}}\left[\min_{\tau_{1}}\left\{\frac{\boldsymbol{r}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}}{2n}\left(\frac{\partial\tilde{\ell}(\tau_{1},0)}{\partial\tau_{1}}\right)^{2}+\tilde{\ell}(\tau_{1},0)\right\}\right]$$

$$-\mathbb{E}_{\boldsymbol{r}=\boldsymbol{\sigma}(\boldsymbol{W}\boldsymbol{x}_{k})}\left[\min_{\tau_{1}}\left\{\frac{\boldsymbol{r}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}}{2n}\left(\frac{\partial\tilde{\ell}(\tau_{1},0)}{\partial\tau_{1}}\right)^{2}+\tilde{\ell}(\tau_{1},0)\right\}\right]$$

$$=o_{d}(1),$$

and therefore from (D.90) we obtain

$$|\mathbb{E}[\Psi_k(\boldsymbol{f}_k) - \Psi_k(\sigma(\boldsymbol{W}\boldsymbol{x}_k))]| \le \frac{o_d(1)}{d},$$

for an absolute constant C > 0, and hence we obtain (D.88).

D.2.7. *Proofs of the Auxiliary Lemmas* Here we provide the proofs of some of the auxiliary lemmas used in our analysis.

Proof of Lemma D.1. We will prove part (b) here, but part (a) will have the exact same proof. Let θ^* be the minimizer of $R_k(\theta, r)$. We can write

$$R_k(\boldsymbol{\theta}^*, \boldsymbol{r}) \leq R_k(\boldsymbol{0}, \boldsymbol{r}),$$

as θ^* is the minimizer.

On the one hand we have

$$R_k(\mathbf{0}, \mathbf{r}) = \frac{1}{n} \sum_{i=1}^n y_i^2 + \|\boldsymbol{\beta}\|_{\ell_2}^2,$$

and thus for any $v \ge 0$:

(D.91)
$$\mathbb{P}\left(R_k(\mathbf{0}, r) \ge v + \|\beta\|_{\ell_2}^2 + 2\mathbb{E}[y_1^2]\right) \le c_1 \exp\left(-nv^2/c_1\right),$$

for an absolute constant $c_1 > 0$.

On the other hand, since $R(\theta, r)$ is λ -strongly convex, and $R(\theta, r) \ge 0$, we can write

$$\|\boldsymbol{\theta}^*\|_{\ell_2}^2 \leq \frac{1}{\lambda} R(\mathbf{0}, \boldsymbol{r}),$$

which together with (D.91) gives use the result.

To prove part (c), we note that r is generated independently from θ_{-k}^* . We can thus write:

$$\mathbb{P}\left(\left|\boldsymbol{r}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*}\right| \geq v\right) \leq \mathbb{P}\left(\left\|\boldsymbol{\theta}_{-k}^{*}\right\|_{\ell_{2}} \geq \sqrt{v}\right) + \mathbb{P}\left(\left|\boldsymbol{r}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*}\right| \geq v\right) \left\|\boldsymbol{\theta}_{-k}^{*}\right\|_{\ell_{2}} < \sqrt{v}\right) \\
\leq c_{2}e^{-v/c_{2}} + \mathbb{P}\left(\left|\boldsymbol{r}^{\mathsf{T}}\boldsymbol{\theta}_{-k}^{*}\right| \geq v\right) \left\|\boldsymbol{\theta}_{-k}^{*}\right\|_{\ell_{2}} < \sqrt{v}\right) \\
\leq c'e^{-v/c'}.$$

where the second step follows from part (a) of the lemma (with c_2 chosen to be sufficiently large); and the last step follows from the independence of r and θ_{-k}^* as well as Lemma D.9.

Also, the proof of part (d) follows simply because

$$\lambda_s \frac{d}{\log(d)} (\mathbf{1}^\mathsf{T} \boldsymbol{\theta}_{-k}^*)^2 \le R_{-k} (\boldsymbol{\theta}_{-k}^*) \le R_{-k} (\mathbf{0}) = \frac{1}{n} \sum_{i \ne k} y_i^2 + \|\boldsymbol{\beta}\|_{\ell_2}^2,$$

where $\mathbf{0}$ is the all-zero vector. By using a bound similar to (D.91) for $R_{-k}(\mathbf{0})$ we obtain the result.

Proof of Lemma D.5. Part (a) is exactly Lemma 12 in [39]. For part (b), consider the matrix \mathbf{R} whose columns are \mathbf{r}_t 's, i.e. $\mathbf{R} = [\mathbf{r}_1 | \mathbf{r}_2 | \cdots | \mathbf{r}_n]$. Since \mathbf{r}_t 's are zero-mean sub-gaussian vectors (see Lemma D.9), we know that its operator norm satisfies:

$$\mathbb{P}\left(||\mathbf{R}|| \ge c_1 \sqrt{d} + v\right) \le c \exp(-v^2/c).$$

Also, define the vector $\alpha = [\alpha_t]_{t \neq k}^T$. Note that $\|\alpha\|_{\ell_2} \leq \sqrt{n}$. We have

$$\mathbb{P}\left(\left\|\frac{1}{n}\sum_{t\neq k}\alpha_{t}\boldsymbol{r}_{t}\right\|_{\ell_{2}}\geq v+c_{1}\right)=\mathbb{P}\left(\left\|\frac{1}{n}\boldsymbol{R}\alpha\right\|_{\ell_{2}}\geq v+c_{1}\right)=\mathbb{P}\left(\left|\left|\boldsymbol{R}\right|\right|\geq v\sqrt{n}+c_{1}\sqrt{n}\right).$$

Given the above relation, and the fact that d and n grow proportionally, the result of the second part of the lemma follows easily.

Part (c) follows from Lemma D.9. Also, part(d) follows from part (c) as well as Lemma D.2 (specifically (D.33)) and the fact that the operator norm of the matrix H_{-k}^{-1} is upper-bounded by $2/\lambda$.

Part (e) follows from the fact that r is a random sub-gaussian vector (see Lemma D.9). We refer to [91] for bounds on the ℓ_2 norm of random sub-gaussian vectors.

To prove part (e), we use (D.33) to write

$$\mathbb{P}\left(\left|\boldsymbol{r}_{t}^{\mathsf{T}}(\tilde{\boldsymbol{\theta}}(\boldsymbol{r}) - \boldsymbol{\theta}_{-k}^{*})\right| \ge \frac{v}{\sqrt{n}}\right) \le \mathbb{P}\left(\frac{1}{n}\left(\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| + \left|\beta_{2}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{p}\right|\right) + \left|\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{e}\right| \ge \frac{v}{\sqrt{n}}\right), \\
(D.92) \le \mathbb{P}\left(\frac{1}{n}\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| \ge \frac{v}{3\sqrt{n}}\right) + \mathbb{P}\left(\frac{1}{n}\left|\beta_{2}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{p}\right| \ge \frac{v}{3\sqrt{n}}\right) + \mathbb{P}\left(\left|\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{e}\right| \ge \frac{v}{3\sqrt{n}}\right)$$

We will now bound each of the terms above. For the first term we have

$$\mathbb{P}\left(\frac{1}{n}\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| \geq \frac{v}{3\sqrt{n}}\right)$$

$$\leq \mathbb{P}\left(\left\|\boldsymbol{r}_{t}\right\|_{\ell_{2}} \geq \sqrt{vn}\right) + \mathbb{P}\left(\frac{1}{n}\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| \geq \frac{v}{3\sqrt{n}}, \left\|\boldsymbol{r}_{t}\right\|_{\ell_{2}} < \sqrt{vn}\right)$$

$$\leq c_{1} \exp\left(-v^{c_{2}}/c_{1}\right) + \mathbb{P}\left(\frac{1}{n}\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| \geq \frac{v}{3\sqrt{n}}, \left\|\boldsymbol{r}_{t}\right\|_{\ell_{2}} < \sqrt{vn}\right),$$

where the last step follows from part (e) and appropriately selecting $c_1, c_2 > 0$. Now, note that the vector \boldsymbol{r} is generated independently from \boldsymbol{r}_t and \boldsymbol{H}_{-k} . As a result, to bound the second term in the RHS of the above relation, we notice that, assuming $\|\boldsymbol{r}_t\|_{\ell_2} < \sqrt{vn}$, we have $\|\beta_1 \boldsymbol{r}_t^\mathsf{T} \boldsymbol{H}_{-k}^{-1}\|_{\ell_2} \le C|\beta_1|\sqrt{vn}$. Hence, we can write

$$\mathbb{P}\left(\frac{1}{n}\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| \geq \frac{v}{3\sqrt{n}} , \|\boldsymbol{r}_{t}\|_{\ell_{2}} < \sqrt{vn}\right) \\
\leq \mathbb{P}\left(\left|\beta_{1}\right| \geq v^{\frac{1}{4}}\right) + \mathbb{P}\left(\frac{1}{n}\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| \geq \frac{v}{3\sqrt{n}} , \|\boldsymbol{r}_{t}\|_{\ell_{2}} < \sqrt{vn} , |\beta_{1}| < v^{\frac{1}{4}}\right) \\
\leq c_{3} \exp(-v^{c_{4}}/c_{3}) + \mathbb{P}\left(\frac{1}{n}\left|\beta_{1}\boldsymbol{r}_{t}^{\mathsf{T}}\boldsymbol{H}_{-k}^{-1}\boldsymbol{r}\right| \geq \frac{v}{3\sqrt{n}} , \|\boldsymbol{r}_{t}\|_{\ell_{2}} < \sqrt{vn} , |\beta_{1}| < v^{\frac{1}{4}}\right) \\
\leq c_{3} \exp(-v^{c_{4}}/c_{3}) + c_{5} \exp(-v^{c_{6}}/c_{5}),$$

where the second inequality follows from (D.34) and by suitably choosing $c_3, c_4 > 0$. The third inequality follows from sub-gaussianity of \boldsymbol{r} , see Lemma D.9, and the fact that, given $|\beta_1| < v^{1/4}$ and $\|\boldsymbol{r}_t\|_{\ell_2} \le \sqrt{vn}$, the random vector $\boldsymbol{u} = \beta_1 \boldsymbol{r}_t \boldsymbol{H}_{-k}^{-1}$ satisfies $\|\boldsymbol{u}\|_{\ell_2} \le C v^{3/4} \sqrt{n}$ and is independently generated from \boldsymbol{r} . Hence, we can bound the second term in the RHS of the above relation by using Lemma D.9 and appropriate choices of $c_5, c_6 > 0$.

The second and third terms in (D.92) can be bounded similarly as the first term but in an easier manner. The second term follows by writing $|\mathbf{r}_t^\mathsf{T} \mathbf{p}| \le ||\mathbf{r}_t||_{\ell_2} ||\mathbf{p}||_{\ell_2}$, and noticing that $||\mathbf{p}||_{\ell_2}$ is upper-bounded by a constant since the norm of \mathbf{J} is bounded and the norm of $\mathbf{\theta}_{-k}^*$ is bounded (see Lemma D.3 and Lemma D.1). Hence, using a similar (but simpler) argument as above, we can write

$$\mathbb{P}\left(\frac{1}{n}\left|\beta_2 \boldsymbol{r}_t^{\mathsf{T}} \boldsymbol{H}_{-k}^{-1} \boldsymbol{p}\right| \ge \frac{v}{3\sqrt{n}}\right) \le c_7 \exp(-v^{c_8}/c_7),$$

for absolute constants c_7 , $c_8 > 8$.

Finally, the third term in the RHS of (D.92) can be bounded by writing $|r_t^\mathsf{T} e| \le ||r_t||_{\ell_2} ||e||_{\ell_2}$ and noticing that $||e||_{\ell_2}$ is small according to (D.34). And, using similar steps as above, we reach to a similar upper bound.

Part (g) follows from Lemma D.1 and Lemma D.6.

Lemma D.9 [Analogous to Lemma 8 in [39]] Assume that $\mathbf{a} = \sigma(\mathbf{W}\mathbf{x})$ and $\mathbf{b} = \mu_1 \mathbf{W}\mathbf{x} + \mu_2 \mathbf{u}$, where \mathbf{x} and \mathbf{u} are generated independently from the normal distribution. Also, let $\mathbf{\Sigma} = \mathbb{E}[\mathbf{b}\mathbf{b}^{\mathsf{T}}]$, i.e. $\mathbf{\Sigma} = \mu_1^2 \mathbf{W} \mathbf{W}^{\mathsf{T}} + \mu_2^2 \mathbf{I}$. Then, there exists an absolute constant c > 0 such that:

(D.93)
$$\mathbb{P}(|\boldsymbol{a}^{\mathsf{T}}\boldsymbol{\beta}| \ge v) \le 2\exp(-\frac{v^2}{c \|\boldsymbol{\beta}\|_{\ell_2}^2 \|\boldsymbol{W}\|^2}),$$

and

(D.94)
$$\mathbb{P}(|\boldsymbol{b}^{\mathsf{T}}\boldsymbol{\beta}| \ge v) \le 2\exp(-\frac{v^2}{c\|\boldsymbol{\beta}\|_{\ell_2}^2\|\boldsymbol{\Sigma}\|}),$$

for a fixed vector $\boldsymbol{\beta} \in \mathbb{R}^d$ and any $v \ge 0$. Here, $\|\boldsymbol{W}\|$ (resp. $\|\boldsymbol{\Sigma}\|$) denotes the operator norm of \boldsymbol{W} (resp. $\boldsymbol{\Sigma}$).

Proof We will be using the following well-known relation: For a *L*-Lipschitz continuous function f and $x \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ we have

(D.95)
$$\mathbb{P}(|f(\boldsymbol{x}) - \mathbb{E}[f(\boldsymbol{x})]| \ge v) \le 2\exp(-\frac{v^2}{4L^2}).$$

Now, note that since σ is the shifted Relu function, it is easy to see that the function $f(x) = \beta^{\mathsf{T}} \sigma(\mathbf{W}x)$ is $\|\beta\|_{\ell_2} \|\mathbf{W}\|$ -Lipschitz continuous. Therefore, we obtain the result using the relation (D.95) and the fact that x is distributed according to the normal distribution.

The proof of the second part can similarly be done by noting that $b = \Sigma^{1/2}\tilde{b}$ where \tilde{b} is distributed according to the standard normal distribution.

Lemma D.10 Let $\mathbf{H} \in \mathbb{R}^{N \times N}$ be such that $||\mathbf{H}|| \le C$ for an absolute constant C > 0. Let $\Sigma_s = \mathbb{E}[\sigma(\mathbf{W}\mathbf{x})\sigma(\mathbf{W}\mathbf{x})^\mathsf{T}]$ and $\Sigma_f = \mathbb{E}[\mathbf{f}\mathbf{f}^\mathsf{T}]$. We have

(D.96)
$$\left| \frac{1}{n} \operatorname{Trace} \left\{ \boldsymbol{H} (\Sigma_s - \Sigma_f) \right\} \right| \leq \frac{c(\log(d))^{3/2}}{\sqrt{d}}.$$

with probability at least $1 - c\exp(-(\log(d))^2/c)$ where c > 0 is an absolute constant.

Proof We first bound each element of the matrix $\Sigma_s - \Sigma_f$. Recall that the k-th element of the vector f is distributed according to

$$\mu_1 \boldsymbol{w}_k^\mathsf{T} \boldsymbol{x} + \mu_2 u_k,$$

where u_k is independently generated from N(0,1) and $\mu_1 = \frac{1}{2}$, $\mu_2 = \sqrt{\frac{1}{4} - \frac{1}{2\pi}}$. As a result, the (ℓ,k) -th element of the matrix Σ_f is

$$\mathbb{E}[(\mu_1 \boldsymbol{w}_k^\mathsf{T} \boldsymbol{x} + \mu_2 u_k)(\mu_1 \boldsymbol{w}_\ell^\mathsf{T} \boldsymbol{x} + \mu_2 u_\ell)] = \mu_1^2 \boldsymbol{w}_k^\mathsf{T} \boldsymbol{w}_\ell + \mu_2^2 \mathbb{I}\{k = \ell\}.$$

Note that $\langle w_\ell, x \rangle$ and $\langle w_k, x \rangle$ are jointly Gaussian with

$$\mathbb{E}[(\boldsymbol{w}_{\ell}^{\mathsf{T}}\boldsymbol{x})^2] = \mathbb{E}[(\boldsymbol{w}_{k}^{\mathsf{T}}\boldsymbol{x})^2] = 1, \quad \mathbb{E}[(\boldsymbol{w}_{\ell}^{\mathsf{T}}\boldsymbol{x})(\boldsymbol{w}_{k}^{\mathsf{T}}\boldsymbol{x})] = \boldsymbol{w}_{k}^{\mathsf{T}}\boldsymbol{w}_{\ell}.$$

Therefore, we have (see e.g., [15, Table 1])

$$\mathbb{E}\left[\sigma(\boldsymbol{w}_{\ell}^{\mathsf{T}}\boldsymbol{x})\sigma(\boldsymbol{w}_{k}^{\mathsf{T}}\boldsymbol{x})\right] = \frac{\sqrt{1 - (\boldsymbol{w}_{k}^{\mathsf{T}}\boldsymbol{w}_{\ell})^{2}} + (\pi - \cos^{-1}(\boldsymbol{w}_{k}^{\mathsf{T}}\boldsymbol{w}_{\ell}))(\boldsymbol{w}_{k}^{\mathsf{T}}\boldsymbol{w}_{\ell})}{2\pi} - \frac{1}{2\pi}$$

$$= \frac{1}{4}\boldsymbol{w}_{\ell}^{\mathsf{T}}\boldsymbol{w}_{k} + (\frac{1}{4} - \frac{1}{2\pi})\mathbb{I}\{k = \ell\} + O((\boldsymbol{w}_{\ell}^{\mathsf{T}}\boldsymbol{w}_{k})^{3})$$

$$= \mu_{1}^{2}\boldsymbol{w}_{\ell}^{\mathsf{T}}\boldsymbol{w}_{k} + \mu_{2}^{2}\mathbb{I}\{k = \ell\} + O\left(\left(\frac{\log(d)}{d}\right)^{\frac{3}{2}}\right),$$
(D.98)

where the last step follows from the fact that with probability at least $1 - c \exp(-(\log(d))^2/c)$ we have for all k, ℓ , such that $k \neq \ell$, we have $|\boldsymbol{w}_{\ell}^{\mathsf{T}} \boldsymbol{w}_{k}| \leq \frac{\log(d)}{d}$. As a result, from (D.97) and (D.98) we obtain

$$\left|\left(\Sigma_s - \Sigma_f\right)_{k,\ell}\right| = O\left(\left(\frac{\log(d)}{d}\right)^{\frac{3}{2}}\right).$$

Hence, the ℓ_2 norm of each column of the matrix $\Sigma_s - \Sigma_f$ is of order $O((\log(d))^{3/2}/d)$. Now, since $||H|| \le C$, then the ℓ_2 norm of each row of H is at most C. Thus, by a simple application of the Cauchy-Schwarz inequality we obtain the result of the lemma.

D.3. Proof of Proposition 6.10 Let us first show that $\widehat{\theta}^*, \widehat{\theta}_{nl}^*$ fall in $\mathcal{C}_{-\theta}$ with high probability for any $\zeta > 0$. We prove the result for $\widehat{\theta}^*$, and remark that the proof is exactly the same for $\widehat{\theta}_{nl}^*$. Consider the objective

$$R(\boldsymbol{\theta}) \coloneqq \frac{1}{2n} \sum_{i=1}^{n} (|y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta}.$$

On the one hand we have

$$R(\mathbf{0}) = \frac{1}{n} \sum_{i=1}^{n} y_i^2$$

where **0** is the all-zero vector. Thus for any $v \ge 0$:

(D.99)
$$\mathbb{P}\left(R_k(\mathbf{0}) \ge v + 2\mathbb{E}[y_1^2]\right) \le c_1 \exp\left(-nv^2/c_1\right),$$

for an absolute constant $c_1 > 0$.

On the other hand, since $R(\theta)$ is ζ -strongly convex, and $R(\theta) \ge 0$, we can write

$$\|\widehat{\boldsymbol{\theta}}^*\|_{\ell_2}^2 \leq \frac{1}{\zeta} R(\mathbf{0}, \boldsymbol{r}),$$

which together with (D.99) proves that $\|\widehat{\boldsymbol{\theta}}^*\|_{\ell_2}$ is bounded above by a constant C with probability at least $1 - e \exp(-cn)$.

To bound $|\mathbf{1}^{\mathsf{T}}\widehat{\boldsymbol{\theta}}^*|$ we note that

$$\zeta \frac{d}{\log(d)} (\mathbf{1}^\mathsf{T} \widehat{\boldsymbol{\theta}}^*)^2 \le R(\widehat{\boldsymbol{\theta}}^*) \le R(\mathbf{0}) = \frac{1}{n} \sum_{i \ne k} y_i^2,$$

By using the bound to (D.99) we obtain with probability $1 - c\exp(-cn)$ that $|\mathbf{1}^T \widehat{\boldsymbol{\theta}}^*| \le C\sqrt{d/(\log(d))}$. Finally, the fact that $\|\widehat{\boldsymbol{\theta}}^*\|_{\ell_\infty}$ is with high probability of order $\sqrt{(\log(d))/d}$ follows in exactly the same manner as the proof of Lemma D.7 and hence we do not repeat the proof here.

We now show that $M(\widehat{\theta}^*) - M(\widehat{\theta}_{nl}^*) \to 0$ in probability. To do so, we use the argument given in [1, Theorem 4]. Assume that $M(\widehat{\theta}^*)$ and $M(\widehat{\theta}_{nl}^*)$ converge to different values, say M_A and M_B . Define $M = (M_A + M_B)/2$ and consider the following optimization problems

$$\bar{\Phi}_A := \min_{\boldsymbol{\theta}: M(\boldsymbol{\theta}) \leq M} \frac{1}{2n} \sum_{i=1}^n (|y_i - \boldsymbol{\theta}^\mathsf{T} \sigma(\boldsymbol{W} \boldsymbol{x}_i)| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \Omega \boldsymbol{\theta},$$

$$\bar{\Phi}_B \coloneqq \min_{\boldsymbol{\theta}: M(\boldsymbol{\theta}) \leq M} \frac{1}{2n} \sum_{i=1}^n (|y_i - \boldsymbol{\theta}^\mathsf{T} \boldsymbol{f}_i| + \varepsilon \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta}.$$

Note that the values $\bar{\Phi}_A$ and $\bar{\Phi}_B$ must be different. Now, using the minimax theorem, and since the above objectives are ζ -strongly convex, we can write

$$\bar{\Phi}_{A} = \sup_{\lambda > 0} -\lambda M + \min_{\boldsymbol{\theta}} \frac{1}{2n} \sum_{i=1}^{n} (|y_{i} - \boldsymbol{\theta}^{\mathsf{T}} \sigma(\boldsymbol{W} \boldsymbol{x}_{i})| + \varepsilon \|\boldsymbol{J} \boldsymbol{\theta}\|_{\ell_{2}})^{2} + \frac{\zeta}{2} \boldsymbol{\theta}^{\mathsf{T}} \Omega \boldsymbol{\theta} + \lambda M(\boldsymbol{\theta}),$$

$$\bar{\Phi}_B = \sup_{\lambda > 0} -\lambda M + \min_{\boldsymbol{\theta}} \frac{1}{2n} \sum_{i=1}^n (|y_i - \boldsymbol{\theta}^\mathsf{T} \boldsymbol{f}_i| + \varepsilon \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{\Omega} \boldsymbol{\theta} + \lambda M(\boldsymbol{\theta}).$$

Now, from the result of Theorem 6.9 we know that for any $\lambda > 0$ the values inside the min converge to the same value. As a result, the quantities $\bar{\Phi}_A$ and $\bar{\Phi}_B$ should converge to the same value (according to [1, Lemma 1]) which is a contradiction with the claim that $M(\widehat{\boldsymbol{\theta}}^*)$ and $M(\widehat{\boldsymbol{\theta}}_{\rm nl}^*)$ converge to different values (i.e. M_A and M_B , respectively). A similar argument can be applied to show that $\|\boldsymbol{J}\widehat{\boldsymbol{\theta}}^*\|_{\ell_2} - \|\boldsymbol{J}\widehat{\boldsymbol{\theta}}_{\rm nl}^*\|_{\ell_2} \to 0$, in probability.

APPENDIX E: PROOFS OF STEP 4: ANALYSIS OF THE GAUSSIAN NOISY LINEAR MODEL VIA CONVEX GAUSSIAN MINIMAX FRAMEWORK

By the Gaussian equivalence property, we henceforth focus on optimization (6.23) and provide a precise characterization of $\overset{\circ}{\mathsf{AR}}_{\mathrm{nl}}(\widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^*)$.

Before proceeding, we will discuss another representation of the model using a few change of variables. Recall from (6.15) that $\boldsymbol{f} \coloneqq \mu_0 \mathbf{1} + \mu_1 \boldsymbol{W} \boldsymbol{x} + \mu_2 \boldsymbol{u}$. For our activation function $\sigma(v) = v \mathbb{I}(v \ge 0) - 1/\sqrt{2\pi}$, we have $\mu_0 = 0$, $\mu_1 = 1/2$ and $\mu_2 = \sqrt{\frac{1}{4} - \frac{1}{2\pi}}$. It is clear that $\boldsymbol{f} \sim \mathsf{N}(0, \boldsymbol{\Sigma})$ with $\boldsymbol{\Sigma} \coloneqq \mu_1^2 \boldsymbol{W} \boldsymbol{W}^\mathsf{T} + \mu_2^2 \boldsymbol{I}$. Also the data generative model (2.1) can be written as:

(E.1)
$$y_i = \langle \mathbf{f}_i, \boldsymbol{\theta}_0 \rangle + w_i, \quad \text{with} \quad w_i \sim \mathsf{N}(0, \sigma^2),$$

for proper choices of σ^2 and θ_0 . Indeed in both models (2.1) and (E.1), $(y_i, f_i) \in \mathbb{R}^{N+1}$ is a centered Gaussian vector. By matching their covariances we obtain

(E.2)
$$\boldsymbol{\Sigma} = \mu_1^2 \boldsymbol{W} \boldsymbol{W}^\mathsf{T} + \mu_2^2 \boldsymbol{I} ,$$

$$\boldsymbol{\theta}_0 = \mu_1 \boldsymbol{\Sigma}^{-1} \boldsymbol{W} \boldsymbol{\beta} ,$$

$$\boldsymbol{\sigma}^2 = \boldsymbol{\tau}^2 + \|\boldsymbol{\beta}\|_{\ell_2}^2 - \mu_1^2 \boldsymbol{\beta}^\mathsf{T} \boldsymbol{W}^\mathsf{T} \boldsymbol{\Sigma}^{-1} \boldsymbol{W} \boldsymbol{\beta} .$$

We next rewrite the objective of optimization (6.23) using this change of variable and also plug in for y_i from (E.1) to obtain

(E.3)
$$\mathcal{L}_{nl}(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} (|\langle \boldsymbol{f}_i, \boldsymbol{\theta}_0 - \boldsymbol{\theta} \rangle + w_i| + \varepsilon ||\boldsymbol{J}\boldsymbol{\theta}||_{\ell_2})^2 + \frac{\zeta}{2} \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{\theta}.$$

We will use a powerful extension of a classical Gaussian process inequality due to Gordon [33] known as *Convex Gaussian Minimax Theorem (CGMT)* [88] to derive a precise asymptotic characterization of $AR_{nl}(\widehat{\theta}_{nl}^*)$. A similar proof technique has been used in [46] to understand the effect of adversarial training on linear regression models. Indeed, for the particular case of $\mu_1 = 0, \mu_2 = 1$ (so $\Sigma = I$) and J = I, the loss function (E.4) reduces to that studied in [46].

The CGMT analysis will output a deterministic scalar optimization which depends on ζ . We need to calculate the solution of this optimization at $\zeta \to 0$. However, as we discuss in our derivation, the objective of this optimization is strongly convex (in minimizing variables) and concave (in maximizing variables). Therefore, by continuity of its solution in the coefficients of the objective, we directly calculate the solution by setting $\zeta = 0$ in the loss $\mathcal{L}_{\rm nl}(\theta)$, bringing us to the following restatement of the loss (with a slight abuse of notation):

(E.4)
$$\mathcal{L}_{\mathrm{nl}}(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left(\left| \langle \boldsymbol{f}_i, \boldsymbol{\theta}_0 - \boldsymbol{\theta} \rangle + w_i \right| + \varepsilon \|\boldsymbol{J}\boldsymbol{\theta}\|_{\ell_2} \right)^2.$$

Consider a change of variable of the form $f_i = \Sigma^{1/2} g_i$ with $g_i \sim N(0, I)$ and $z = \Sigma^{1/2} (\theta - \theta_0)$. Also define

$$\ell(v; \boldsymbol{\theta}) \coloneqq \frac{1}{2} (|v| + \varepsilon ||\boldsymbol{J}\boldsymbol{\theta}||_{\ell_2})^2.$$

Then, optimization problem (E.4) can be equivalently written in the form

(E.5)
$$\min_{\boldsymbol{z} \in \mathbb{R}^{N}, \boldsymbol{v} \in \mathbb{R}^{n}} \frac{1}{n} \sum_{i=1}^{n} \ell\left(v_{i}; \boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z}\right) \text{ subject to } \boldsymbol{v} = \boldsymbol{w} - \boldsymbol{G} \boldsymbol{z}.$$

By writing the dual of this optimization problem (with dual variable $\frac{u}{\sqrt{d}}$) we get

(E.6)
$$\min_{\boldsymbol{z} \in \mathbb{R}^{N}, \boldsymbol{v} \in \mathbb{R}^{n}} \max_{\boldsymbol{u} \in \mathbb{R}^{n}} \frac{1}{\sqrt{d}} \left\{ \boldsymbol{u}^{\mathsf{T}} \boldsymbol{G} \boldsymbol{z} - \boldsymbol{u}^{T} \boldsymbol{w} + \boldsymbol{u}^{T} \boldsymbol{v} \right\} + \frac{1}{n} \sum_{i=1}^{n} \ell \left(v_{i}; \boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z} \right)$$

$$= \min_{\boldsymbol{z} \in \mathbb{R}^{N}, \boldsymbol{v} \in \mathbb{R}^{n}} \max_{\boldsymbol{u} \in \mathbb{R}^{n}} \frac{1}{\sqrt{d}} \left\{ \boldsymbol{u}^{\mathsf{T}} \boldsymbol{G} \boldsymbol{z} - \boldsymbol{u}^{\mathsf{T}} \boldsymbol{w} + \boldsymbol{u}^{\mathsf{T}} \boldsymbol{v} \right\} + \bar{\ell}(\boldsymbol{v}; \boldsymbol{z}),$$

where

$$\bar{\ell}(\boldsymbol{v}; \boldsymbol{z}) \coloneqq \frac{1}{n} \sum_{i=1}^{n} \ell\left(v_{i}; \boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z}\right) \\
= \frac{1}{2n} \left\|\boldsymbol{v}\right\|_{\ell_{2}}^{2} + \frac{\varepsilon}{n} \left\|\boldsymbol{v}\right\|_{\ell_{1}} \left\|\boldsymbol{J}(\boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z})\right\|_{\ell_{2}} + \frac{\varepsilon^{2}}{2} \left\|\boldsymbol{J}(\boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z})\right\|_{\ell_{2}}^{2}.$$

The minimax optimization (E.6) is in a form that we can apply the CGMT framework. Formally, the CGMT framework concerns problems of the form

(E.7)
$$\min_{\boldsymbol{z} \in \mathcal{S}_{\boldsymbol{z}}} \max_{\boldsymbol{u} \in \mathcal{S}_{\boldsymbol{u}}} \boldsymbol{u}^T \boldsymbol{G} \boldsymbol{z} + \psi(\boldsymbol{z}, \boldsymbol{u}),$$

with G a matrix with i.i.d standard normal entries and shows that this problem is asymptotically equivalent to the following problem:

(E.8)
$$\min_{\boldsymbol{z} \in \mathcal{S}_{\boldsymbol{z}}} \max_{\boldsymbol{u} \in \mathcal{S}_{\boldsymbol{u}}} \|\boldsymbol{z}\|_{\ell_2} \boldsymbol{g}^T \boldsymbol{u} + \|\boldsymbol{u}\|_{\ell_2} \boldsymbol{h}^T \boldsymbol{z} + \psi(\boldsymbol{z}, \boldsymbol{u}),$$

where g and h are independent Gaussian vectors with i.i.d. $\mathcal{N}(0,1)$ entries and $\psi(z,u)$ is convex in z and concave in u. Here, the sets \mathcal{S}_z and \mathcal{S}_u are compact sets. We refer to [88, Theorem 3] for precise statements regarding the equivalence of (E.7) and (E.8).

Following [88] we shall refer to problems of the form (E.7) as the *Primal Problem (PO)* and refer to problems of the form (E.8) as the *Auxiliary Problem (AO)*.

As described above the CGMT framework requires the minimization/maximization to be over compact sets. This technical issue can be avoided by a common trick in this literature where one introduces "artificial" boundedness constraint which do not effect the optimal solution. Specifically, following [88] we can add constraints of the form $S_z = \{z \mid \|z\|_{\ell_2} \le K_\alpha\}$ and $S_u = \{u : \|u\|_{\ell_2} \le K_\beta\}$ for sufficiently large constants K_α and K_β without changing the optimal solution of (E.6) in a precise asymptotic sense. We leave out a detailed argument here and refer to [46, Appendix B] for similar arguments. This allows us to replace (E.6) with

(E.9)
$$\min_{\boldsymbol{z} \in S_{\boldsymbol{z}}, \boldsymbol{v} \in \mathbb{R}^n} \max_{\boldsymbol{u} \in S_{\boldsymbol{u}}} \frac{1}{\sqrt{d}} \left\{ \boldsymbol{u}^\mathsf{T} \boldsymbol{G} \boldsymbol{z} - \boldsymbol{u}^\mathsf{T} \boldsymbol{w} + \boldsymbol{u}^\mathsf{T} \boldsymbol{v} \right\} + \bar{\ell}(\boldsymbol{v}; \boldsymbol{z}).$$

Observe that the above loss function has a bilinear term $u^T G z$, with $G_{ij} \sim N(0,1)$ independently, plus a function of the form

$$\psi(\boldsymbol{z}, \boldsymbol{v}, \boldsymbol{u}) \coloneqq \frac{1}{\sqrt{d}} \left\{ -\boldsymbol{u}^\mathsf{T} \boldsymbol{w} + \boldsymbol{u}^\mathsf{T} \boldsymbol{v} \right\} + \bar{\ell}(\boldsymbol{v}; \boldsymbol{z}),$$

which is jointly convex in (z, v) and concave in u.

Therefore the corresponding AO problem takes the following form

(E.10)
$$\min_{\boldsymbol{z} \in S_{\boldsymbol{z}}, \boldsymbol{v}} \max_{\boldsymbol{u} \in S_{\boldsymbol{u}}} \frac{1}{\sqrt{d}} \left\{ \|\boldsymbol{z}\|_{\ell_2} \boldsymbol{g}^T \boldsymbol{u} + \|\boldsymbol{u}\|_{\ell_2} \boldsymbol{h}^T \boldsymbol{z} - \boldsymbol{u}^\mathsf{T} \boldsymbol{w} + \boldsymbol{u}^\mathsf{T} \boldsymbol{v} \right\} + \bar{\ell}(\boldsymbol{v}; \boldsymbol{z}).$$

This concludes the derivation of the AO problem.

E.1. Scalarization of the AO problem We next simplify the AO problem by considering this problem in the asymptotic regime. We start by maximizing over u. Write $u = \beta \widetilde{u}$ with $\widetilde{u} \in \mathbb{S}^{n-1}$ and $0 \le \beta \le K_{\beta}$. Using this decomposition we have

$$\begin{aligned} & \max_{\boldsymbol{u} \in \mathcal{S}_{\boldsymbol{u}}} \|\boldsymbol{z}\|_{\ell_{2}} \boldsymbol{g}^{T} \boldsymbol{u} + \|\boldsymbol{u}\|_{\ell_{2}} \boldsymbol{h}^{T} \boldsymbol{z} - \boldsymbol{u}^{T} \boldsymbol{w} + \boldsymbol{u}^{T} \boldsymbol{v} \\ &= \max_{0 \leq \beta \leq K_{\beta}} \max_{\widetilde{\boldsymbol{u}} \in \mathbb{S}^{n-1}} \beta \|\boldsymbol{z}\|_{\ell_{2}} \boldsymbol{g}^{T} \widetilde{\boldsymbol{u}} + \beta \boldsymbol{h}^{T} \boldsymbol{z} - \beta \widetilde{\boldsymbol{u}}^{T} \boldsymbol{w} + \beta \widetilde{\boldsymbol{u}}^{T} \boldsymbol{v} \\ &= \max_{0 \leq \beta \leq K_{\beta}} \max_{\widetilde{\boldsymbol{u}} \in \mathbb{S}^{n-1}} \beta \widetilde{\boldsymbol{u}}^{T} (\|\boldsymbol{z}\|_{\ell_{2}} \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}) + \beta \boldsymbol{h}^{T} \boldsymbol{z} \\ &= \max_{0 \leq \beta \leq K_{\beta}} \beta \|\|\boldsymbol{z}\|_{\ell_{2}} \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}} + \beta \boldsymbol{h}^{T} \boldsymbol{z}. \end{aligned}$$

After substituting the above into AO problem (E.10), it reads

$$\min_{oldsymbol{z} \in \mathcal{S}_{oldsymbol{z}}, oldsymbol{v}} \max_{0 \leq eta \leq K_{eta}} \frac{eta}{\sqrt{d}} \| \| oldsymbol{z} \|_{\ell_2} oldsymbol{g} - oldsymbol{w} + oldsymbol{v} \|_{\ell_2} + \frac{eta}{\sqrt{d}} oldsymbol{h}^T oldsymbol{z} + ar{\ell}(oldsymbol{v}; oldsymbol{z}).$$

We next aim to simplify the minimization over v and z, but a hurdle is that they are coupled through the term $\ell(v;z)$. To address this technical issue, we consider the conjugate of $\bar{\ell}(v;z)$ in with respect to z. That is,

$$\bar{\ell}(\boldsymbol{v};\boldsymbol{z}) = \sup_{\boldsymbol{q}} \boldsymbol{q}^T \boldsymbol{z} - \widetilde{\ell}(\boldsymbol{v};\boldsymbol{q})$$

The AO problem then can be written as

(E.11)
$$\min_{\boldsymbol{z} \in \mathcal{S}_{\boldsymbol{z}}, \boldsymbol{v}} \max_{0 \le \beta \le K_{\beta}, \boldsymbol{q}} \frac{\beta}{\sqrt{d}} \| \| \boldsymbol{z} \|_{\ell_{2}} \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v} \|_{\ell_{2}} + \frac{\beta}{\sqrt{d}} \boldsymbol{h}^{T} \boldsymbol{z} + \boldsymbol{q}^{T} \boldsymbol{z} - \widetilde{\ell}(\boldsymbol{v}; \boldsymbol{q}).$$

In the above optimization the order of minimization and maximization can be flipped using the Sion's theorem and the fact that the original PO problem is convex/concave in the min/max parameters. The argument only uses the convexity of the loss $\ell(v;q)$ and we refer to [86, Appendix A.2.4] for a detailed argument. This brings us to

$$\max_{0 \leq \beta \leq K_{\beta}, \boldsymbol{q}} \min_{\boldsymbol{z} \in \mathcal{S}_{\boldsymbol{z}}, \boldsymbol{v}} \frac{\beta}{\sqrt{d}} \| \| \boldsymbol{z} \|_{\ell_2} \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v} \|_{\ell_2} + \frac{\beta}{\sqrt{d}} \boldsymbol{h}^T \boldsymbol{z} + \boldsymbol{q}^T \boldsymbol{z} - \widetilde{\ell}(\boldsymbol{v}; \boldsymbol{q}).$$

We optimize over the direction and norm of z ($||z||_{\ell_2} = \alpha$) to get

(E.12)
$$\max_{0 \le \beta \le K_{\beta}, \mathbf{q}} \min_{0 \le \alpha \le K_{\alpha}, \mathbf{v}} \frac{\beta}{\sqrt{d}} \|\alpha \mathbf{g} - \mathbf{w} + \mathbf{v}\|_{\ell_{2}} - \alpha \left\| \frac{\beta}{\sqrt{d}} \mathbf{h} + \mathbf{q} \right\|_{\ell_{\alpha}} - \widetilde{\ell}(\mathbf{v}; \mathbf{q}).$$

Note that $\widetilde{\ell}(v;q)$ is convex in q and so the AO objective (E.12) is clearly jointly concave in q and β . Also since $\overline{\ell}$ is jointly convex in (v,z), then $-\overline{\ell}(v;z)$ is jointly concave in (v,z). Also q^Tz is jointly concave in (v,z). Therefore, $q^Tz - \overline{\ell}(v;z)$ is jointly concave in (v,z) and based on the partial maximization rule we can conclude that $\widetilde{\ell}(v;q)$ should be concave in v. The other terms are also trivially jointly convex in α, v so that overall the objective is jointly convex in α, v . Therefore, by virtue of Sion's min-max Theorem [76]) we can change the order of the mins and maxs as we please and rewrite the AO problem as

$$\min_{0 \le \alpha \le K_{\alpha}, \boldsymbol{v}} \max_{0 \le \beta \le K_{\beta}, \boldsymbol{q}} \frac{\beta}{\sqrt{d}} \|\alpha \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}} - \alpha \left\| \frac{\beta}{\sqrt{d}} \boldsymbol{h} + \boldsymbol{q} \right\|_{\ell_{2}} - \widetilde{\ell}(\boldsymbol{v}; \boldsymbol{q}).$$

To continue we shall calculate the conjugate function $\widetilde{\ell}$. This is the subject of the next lemma.

Lemma E.1 The conjugate of

$$\bar{\ell}(\boldsymbol{v};\boldsymbol{z}) \coloneqq \frac{1}{2n} \sum_{i=1}^{n} \left(|v_i| + \varepsilon \left\| \boldsymbol{J} (\boldsymbol{\theta}_0 + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z}) \right\|_{\ell_2} \right)^2,$$

with respect to the variable z is given by

$$\widetilde{\ell}(\boldsymbol{v};\boldsymbol{q}) \coloneqq \sup_{\boldsymbol{z}} \boldsymbol{q}^T \boldsymbol{z} - \overline{\ell}(\boldsymbol{v};\boldsymbol{z}) = -\langle \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0, \boldsymbol{q} \rangle + \frac{1}{2} \left(\frac{1}{\varepsilon} \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{J}^{-1} \boldsymbol{q} \right\|_{\ell_2} - \frac{1}{n} \left\| \boldsymbol{v} \right\|_{\ell_1} \right)_+^2 - \frac{1}{2n} \left\| \boldsymbol{v} \right\|_{\ell_2}^2.$$

We refer to Section E.5 for the proof of this lemma. Plugging in for $\widetilde{\ell}$ from Lemma E.1 in the AO problem we arrive at

$$\min_{0 \le \alpha < K_{\alpha}, \boldsymbol{v}} \max_{0 \le \beta \le K_{\beta}, \boldsymbol{q}} \frac{\beta}{\sqrt{d}} \|\alpha \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}} - \alpha \left\| \frac{\beta}{\sqrt{d}} \boldsymbol{h} + \boldsymbol{q} \right\|_{\ell_{2}} \\
+ \langle \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0}, \boldsymbol{q} \rangle - \frac{1}{2} \left(\frac{1}{\varepsilon} \|\boldsymbol{J}^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{q} \|_{\ell_{2}} - \frac{\|\boldsymbol{v}\|_{\ell_{1}}}{n} \right)_{+}^{2} + \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_{2}}^{2}.$$
(E.13)

Optimization over q: To simplify the AO problem further, we next focus on maximization over q. Consider the change of variable $\tilde{q} := J^{-1} \Sigma^{1/2} q$ and keep only the terms in the AO objective which involve q.

$$\begin{aligned} & \max_{\boldsymbol{q}} \quad -\alpha \left\| \frac{\beta}{\sqrt{d}} \boldsymbol{h} + \boldsymbol{q} \right\|_{\ell_{2}} + \left\langle \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0}, \boldsymbol{q} \right\rangle - \frac{1}{2} \left(\frac{1}{\varepsilon} \left\| \boldsymbol{J}^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{q} \right\|_{\ell_{2}} - \frac{\| \boldsymbol{v} \|_{\ell_{1}}}{n} \right)_{+}^{2} \\ & = \max_{\boldsymbol{q}} \quad -\alpha \left\| \frac{\beta}{\sqrt{d}} \boldsymbol{h} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{J} \boldsymbol{\widetilde{q}} \right\|_{\ell_{2}} + \left\langle \boldsymbol{\theta}_{0}, \boldsymbol{J} \boldsymbol{\widetilde{q}} \right\rangle - \frac{1}{2} \left(\frac{1}{\varepsilon} \left\| \boldsymbol{\widetilde{q}} \right\|_{\ell_{2}} - \frac{\| \boldsymbol{v} \|_{\ell_{1}}}{n} \right)_{+}^{2} \\ & = \max_{\boldsymbol{\widetilde{q}}, 0 \le \tau_{q}} \quad -\frac{\alpha}{2\tau_{q}} \left\| \frac{\beta}{\sqrt{d}} \boldsymbol{h} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{J} \boldsymbol{\widetilde{q}} \right\|_{\ell_{2}}^{2} - \frac{\alpha\tau_{q}}{2} + \left\langle \boldsymbol{\theta}_{0}, \boldsymbol{J} \boldsymbol{\widetilde{q}} \right\rangle - \frac{1}{2} \left(\frac{1}{\varepsilon} \left\| \boldsymbol{\widetilde{q}} \right\|_{\ell_{2}} - \frac{\| \boldsymbol{v} \|_{\ell_{1}}}{n} \right)_{+}^{2} \\ & = \max_{\boldsymbol{\widetilde{q}}, 0 \le \tau_{q}} \quad -\frac{\alpha}{2\tau_{q}} \left\| \frac{\beta}{\sqrt{d}} \boldsymbol{h} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{J} \boldsymbol{\widetilde{q}} - \frac{\tau_{q}}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} \right\|_{\ell_{2}}^{2} + \frac{\tau_{q}}{2\alpha} \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} \right\|_{\ell_{2}}^{2} - \left\langle \frac{\beta}{\sqrt{d}} \boldsymbol{h}, \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} \right\rangle \\ & - \frac{\alpha\tau_{q}}{2} - \frac{1}{2} \left(\frac{1}{\varepsilon} \left\| \boldsymbol{\widetilde{q}} \right\|_{\ell_{2}} - \frac{\| \boldsymbol{v} \|_{\ell_{1}}}{n} \right)_{+}^{2} \end{aligned}$$

We next maximize over \widetilde{q} by introducing a new dummy variable γ for $\|\widetilde{q}\|_{\ell_2}$. This brings us to the following problem

(E.14)
$$\min_{\widetilde{\boldsymbol{q}},0 \le \gamma} \frac{\alpha}{2\tau_q} \left\| \frac{\beta}{\sqrt{d}} \boldsymbol{h} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{J} \widetilde{\boldsymbol{q}} - \frac{\tau_q}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 \right\|_{\ell_2}^2 + \frac{1}{2} \left(\frac{1}{\varepsilon} \gamma - \frac{\|\boldsymbol{v}\|_{\ell_1}}{n} \right)_+^2$$
subject to $\|\widetilde{\boldsymbol{q}}\|_{\ell_2} = \gamma$.

We continue by the following lemma and refer to Section E.5 for its proof.

Lemma E.2 Let $H \in \mathbb{R}^{d \times d}$ be invertible and $\mathbf{r} \in \mathbb{R}^d$, $c_0, c_1 \in \mathbb{R}$. Consider the following optimization problem:

(E.15)
$$\min_{\widetilde{\boldsymbol{q}},0 \leq \gamma} \frac{c_0}{2} \|\boldsymbol{H}\widetilde{\boldsymbol{q}} - \boldsymbol{r}\|_{\ell_2}^2 + \frac{1}{2} \left(\frac{1}{\varepsilon} \gamma - c_1\right)_+^2$$
s.t. $\|\widetilde{\boldsymbol{q}}\|_{\ell_2} = \gamma$

Define

(E.16)
$$Q(\boldsymbol{H}, \boldsymbol{r}, \gamma) = \sup_{\lambda > 0} \frac{\lambda}{2} (\boldsymbol{r}^{\mathsf{T}} (\boldsymbol{H} \boldsymbol{H}^{\mathsf{T}} + \lambda \boldsymbol{I})^{-1} \boldsymbol{r} - \gamma^{2}).$$

Then, the optimal objective value of (E.15) is given by

$$\min_{\gamma \geq 0} c_0 Q(\boldsymbol{H}, \boldsymbol{r}, \gamma) + \frac{1}{2} \left(\frac{1}{\varepsilon} \gamma - c_1 \right)_+^2.$$

Using Lemma E.2, the optimal value of (E.14) is given by

$$\min_{\gamma \geq 0} \frac{\alpha}{\tau_q} Q(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}, \frac{\tau_q}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \boldsymbol{h}, \gamma) + \frac{1}{2} \left(\frac{\gamma}{\varepsilon} - \frac{\|\boldsymbol{v}\|_{\ell_1}}{n} \right)^2.$$

Therefore, by substituting in (E.12) the AO optimization can be simplified as

$$\min_{0 \le \alpha < K_{\alpha}, \boldsymbol{v}} \max_{0 \le \beta \le K_{\beta}, 0 \le \gamma, \tau_{q}} \frac{\beta}{\sqrt{d}} \|\alpha \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}} - \frac{\alpha}{\tau_{q}} Q(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}, \frac{\tau_{q}}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} - \frac{\beta}{\sqrt{d}} \boldsymbol{h}, \gamma) \\
+ \frac{\tau_{q}}{2\alpha} \|\boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0}\|_{\ell_{2}}^{2} - \langle \frac{\beta}{\sqrt{d}} \boldsymbol{h}, \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} \rangle - \frac{\alpha \tau_{q}}{2} - \frac{1}{2} (\frac{\gamma}{\varepsilon} - \frac{\|\boldsymbol{v}\|_{\ell_{1}}}{n})_{+}^{2} + \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_{2}}^{2}.$$

Before proceeding further with our simplification of the AO problem, let us state the following lemma which is used to discuss the convexity-concavity of the objective and justification of changing the order of maximization and minimization. We refer to Section E.5 for its proof.

Lemma E.3 The function

$$f(\gamma, \beta, \tau_q) \coloneqq \frac{1}{\tau_q} Q(\mathbf{\Sigma}^{-1/2} \mathbf{J}, \frac{\tau_q}{\alpha} \mathbf{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \mathbf{h}, \gamma),$$

is jointly convex in the variables $(\gamma, \frac{\beta}{\sqrt{d}}, \tau_q)$.

As a result this lemma, the objective (E.17) is jointly concave in (γ, β, τ_q) . Also recall that since $\widetilde{\ell}$ was concave the objective (E.12) was jointly convex in (α, v) . Since maximization (with respect to direction of \widetilde{q}) preserves convexity (pointwise maximum of convex functions is convex), therefore the objective (E.17) is jointly convex in (α, v) .

Therefore, by another use of Sion's min-max theorem, we can change the order of min and max in (E.17) and write it equivalently as

$$\max_{0 \leq \beta \leq K_{\beta}, 0 \leq \gamma, \tau_{q}} \min_{0 \leq \alpha < K_{\alpha}, \boldsymbol{v}} \frac{\beta}{\sqrt{d}} \|\alpha \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}} - \frac{\alpha}{\tau_{q}} Q(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}, \frac{\tau_{q}}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} - \frac{\beta}{\sqrt{d}} \boldsymbol{h}, \gamma) \\
(E.18) + \frac{\tau_{q}}{2\alpha} \|\boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0}\|_{\ell_{2}}^{2} - \frac{1}{\sqrt{d}} \langle \beta \boldsymbol{h}, \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} \rangle - \frac{\alpha \tau_{q}}{2} - \frac{1}{2} \left(\frac{\gamma}{\varepsilon} - \frac{\|\boldsymbol{v}\|_{\ell_{1}}}{n}\right)_{+}^{2} + \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_{2}}^{2}.$$

We next focus on minimization over v. Keeping only the terms in (E.17) that depend on v we have

$$\min_{\boldsymbol{v}} \frac{\beta}{\sqrt{d}} \|\alpha \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}} + \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_{2}}^{2} - \frac{1}{2} \left(\frac{\gamma}{\varepsilon} - \frac{\|\boldsymbol{v}\|_{\ell_{1}}}{n}\right)_{+}^{2}$$
(E.19)
$$= \min_{\tau_{g} \geq 0, \boldsymbol{v}} \frac{\beta}{2n\tau_{g}} \|\alpha \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}}^{2} + \frac{\beta\tau_{g}n}{2d} + \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_{2}}^{2} - \frac{1}{2n^{2}} \left(\frac{n\gamma}{\varepsilon} - \|\boldsymbol{v}\|_{\ell_{1}}\right)_{+}^{2}.$$

Recall the definition of the Moreau envelope function of a function f at a point x with parameter μ ,

$$e_f(\boldsymbol{x}; \rho) \equiv \min_{\boldsymbol{v}} \frac{1}{2\rho} \|\boldsymbol{x} - \boldsymbol{v}\|_{\ell_2}^2 + f(\boldsymbol{v}).$$

and define

(E.20)
$$f(\boldsymbol{v};\gamma) \coloneqq \frac{1}{2} \|\boldsymbol{v}\|_{\ell_2}^2 - \frac{1}{2n} (\frac{n}{\varepsilon} \gamma - \|\boldsymbol{v}\|_{\ell_1})_+^2.$$

Note that $f(v; \gamma)$ is convex in v (since $-\ell(v; q)$ was convex in v). Thus, (E.19) can be rewritten in the more compact form

(E.21)
$$\min_{\tau_g \ge 0} \frac{1}{n} e_f \left(\boldsymbol{w} - \alpha \boldsymbol{g}; \frac{\tau_g}{\beta} \right) + \frac{\beta \tau_g}{2} \frac{n}{d}.$$

We next invoke the result of [46, Lemma 6.3] which gives a characterization of the Moreau envelope function $e_f(x; \mu)$.

Lemma E.4 ([46, Lemma 6.3]) Consider the function f given by (E.20). Then,

$$e_f(x; \rho) = \frac{1}{2(\rho+1)} \|x\|_{\ell_2}^2 + \min_{\nu \ge 0} G_n(x; \rho, \gamma, \nu),$$

where

$$(E.22) \quad G_n(\boldsymbol{x}; \rho, \gamma, \nu) = \frac{1}{2\rho(\rho+1)} \|\boldsymbol{x} - \mathsf{ST}(\boldsymbol{x}; \nu)\|_{\ell_2}^2 - \frac{1}{2n} \left(\frac{n}{\varepsilon} \gamma - \frac{1}{1+\rho} \|\mathsf{ST}(\boldsymbol{x}; \nu)\|_{\ell_1} \right)_{\perp}^2,$$

and $ST(x; \nu)$ is the soft-thresholding function defined as

$$[\mathsf{ST}(\boldsymbol{x};\nu)]_i = \begin{cases} x_i - \lambda, & \text{if} & x_i \ge \lambda, \\ 0 & \text{if} & |x_i| \le \lambda, \\ x_i + \lambda & \text{if} & x_i \le -\lambda, \end{cases}$$

for each coordinate i. Furthermore, $e_f(x;\tau)$ is strictly convex in x.

Using this characterization in (E.19) we get

$$\min_{\boldsymbol{v}} \frac{\beta}{\sqrt{d}} \|\alpha \boldsymbol{g} - \boldsymbol{w} + \boldsymbol{v}\|_{\ell_{2}} + \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_{2}}^{2} - \frac{1}{2} \left(\frac{\gamma}{\varepsilon} - \frac{\|\boldsymbol{v}\|_{\ell_{1}}}{n}\right)_{+}^{2}$$

$$= \min_{\tau_{g} \geq 0} \frac{1}{n} e_{f} \left(\boldsymbol{w} - \alpha \boldsymbol{g}; \frac{\tau_{g}}{\beta}\right) + \frac{\beta \tau_{g}}{2} \frac{n}{d}$$

$$= \min_{\tau_{g} \geq 0} \frac{\beta \tau_{g}}{2} \frac{n}{d} + \frac{1}{n} \frac{\beta}{2(\tau_{g} + \beta)} \|\boldsymbol{w} - \alpha \boldsymbol{g}\|_{\ell_{2}}^{2} + \frac{1}{n} \min_{\nu \geq 0} G_{n} (\boldsymbol{w} - \alpha \boldsymbol{g}; \frac{\tau_{g}}{\beta}, \gamma, \nu).$$
(E.23)

Next by plugging (E.23) in (E.18) we arrive at the following AO formulation:

$$\max_{0 \leq \beta \leq K_{\beta}, 0 \leq \gamma, \tau_{q}} \min_{0 \leq \alpha < K_{\alpha}, 0 \leq \tau_{g}, \nu} -\frac{\alpha}{\tau_{q}} Q(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}, \frac{\tau_{q}}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} - \frac{\beta}{\sqrt{d}} \boldsymbol{h}, \gamma) + \frac{\tau_{q}}{2\alpha} \|\boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0}\|_{\ell_{2}}^{2} - \frac{\alpha \tau_{q}}{2} \right) \\
(E.24) -\frac{1}{\sqrt{d}} \langle \beta \boldsymbol{h}, \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} \rangle + \frac{\beta \tau_{g}}{2} \frac{n}{d} + \frac{\beta}{2(\tau_{g} + \beta)} \frac{1}{n} \|\boldsymbol{w} - \alpha \boldsymbol{g}\|_{\ell_{2}}^{2} \\
+ \frac{1}{n} G_{n} (\boldsymbol{w} - \alpha \boldsymbol{g}; \frac{\tau_{g}}{\beta}, \gamma, \nu).$$

Recall that the problem (E.19) was jointly convex in (v, α, τ_g) and (E.18) jointly concave in (β, γ, τ_q) . Since partial minimization preserves convexity we therefore conclude that the objective (E.24) is jointly convex in (α, τ_g) and jointly concave in (β, γ, τ_q) (after the minimization over $\nu \ge 0$ has been carried out).

E.2. Convergence analysis of the AO problem

E.2.1. *Pointwise convergence* We next derive the pointwise limit of the AO objective in the asymptotic regime that $N/d \to \psi_1$ and $n/d \to \psi_2$, as $n \to \infty$.

Recalling the definition of θ_0 from (E.2) we have

$$\left\|\boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_{0}\right\|_{\ell_{2}}^{2} = \mu_{1}^{2} \left\|\boldsymbol{\Sigma}^{-1/2}\boldsymbol{W}\boldsymbol{\beta}\right\|_{\ell_{2}}^{2} = \frac{\mu_{1}^{2} \left\|\boldsymbol{\beta}\right\|_{\ell_{2}}^{2}}{d} \operatorname{trace}(\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\Sigma}\boldsymbol{W}),$$

where we used the fact that the distribution of W is rotationally invariant. By our assumption $\|\beta\|_{\ell_2} \to 1$. Let $0 \le s_1, \dots, s_N$ denote the eigenvalues of WW^T . By invoking the definition of Σ from (E.2) we have

$$\begin{split} \left\| \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} \right\|_{\ell_{2}}^{2} &= \frac{\psi_{1} \mu_{1}^{2} \|\boldsymbol{\beta}\|_{\ell_{2}}^{2}}{N} \operatorname{trace}(\boldsymbol{W}^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{W}) \\ &= \frac{\psi_{1} \mu_{1}^{2} \|\boldsymbol{\beta}\|_{\ell_{2}}^{2}}{N} \sum_{i=1}^{d} \frac{s_{i}}{\mu_{1}^{2} s_{i} + \mu_{2}^{2}} \\ (E.26) &= \frac{\psi_{1} \mu_{1}^{2} \|\boldsymbol{\beta}\|_{\ell_{2}}^{2}}{N} \sum_{i=1}^{d} \frac{1}{\mu_{1}^{2}} \left(1 - \frac{\mu_{2}^{2} / \mu_{1}^{2}}{s_{i} + \mu_{2}^{2} / \mu_{1}^{2}} \right) \rightarrow \psi_{1} \left(1 + \frac{\mu_{2}^{2}}{\mu_{1}^{2}} S \left(- \frac{\mu_{2}^{2}}{\mu_{1}^{2}} ; \psi_{1} \right) \right), \end{split}$$

in probability where $S(z) = \int \frac{\rho(s)}{z-s} ds$ is the Stieltjes transform of the spectral density ρ of the matrix WW^T . The formula for S(z) is given in Proposition F.2 and since it is a function of ψ_1 , we make this dependence clear in the notation and write $S(z; \psi_1)$ henceforth.

Since for our activation $\mu_1 = \frac{1}{2}$ and $\mu_2 = \sqrt{\frac{1}{4} - \frac{1}{2\pi}}$, this simplifies to

$$\left\| \mathbf{\Sigma}^{1/2} \boldsymbol{\theta}_0 \right\|_{\ell_2}^2 \to \psi_1 \left(1 + \left(1 - \frac{2}{\pi} \right) S \left(\frac{2}{\pi} - 1; \psi_1 \right) \right),$$

in probability. This together with (E.2) implies that

(E.28)
$$\sigma^{2} = \tau^{2} + \|\boldsymbol{\beta}\|_{\ell_{2}}^{2} - \|\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\theta}_{0}\|_{\ell_{2}}^{2} \to \tau^{2} + 1 - \psi_{1}\left(1 + \left(1 - \frac{2}{\pi}\right)S\left(\frac{2}{\pi} - 1; \psi_{1}\right)\right).$$

We next note that since $\|\mathbf{\Sigma}^{1/2}\boldsymbol{\theta}_0\|_{\ell_2} = O(1)$, for $\boldsymbol{h} \sim \mathsf{N}(0, \boldsymbol{I})$ we have

(E.29)
$$\frac{1}{\sqrt{n}}\langle \boldsymbol{h}, \boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_0 \rangle \to 0,$$

in probability, as $n \to \infty$. In addition, since $\mathbf{g} \sim \mathsf{N}(0, \mathbf{I}_n)$ and $\mathbf{w} \sim \mathsf{N}(0, \sigma^2 \mathbf{I}_n)$, we have

(E.30)
$$\frac{1}{n} \|\alpha \boldsymbol{g} - \boldsymbol{w}\|_{\ell_2}^2 \to \alpha^2 + \sigma^2,$$

in probability.

We next proceed by calculating the limit of the Q function. By definition,

$$Q(\boldsymbol{\Sigma}^{-1/2}\boldsymbol{J}, \frac{\tau_q}{\alpha}\boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}}\boldsymbol{h}, \gamma)$$

$$= \sup_{\lambda \geq 0} \frac{\lambda}{2} \left[(\frac{\tau_q}{\alpha}\boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}}\boldsymbol{h})^{\mathsf{T}} (\boldsymbol{\Sigma}^{-1/2}\boldsymbol{J}^2\boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I})^{-1} (\frac{\tau_q}{\alpha}\boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}}\boldsymbol{h}) - \gamma^2 \right]$$
(E.31)
$$= \sup_{\lambda \geq 0} \frac{\lambda}{2} \left[(\frac{\tau_q}{\alpha}\boldsymbol{\Sigma}\boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}}\boldsymbol{\Sigma}^{1/2}\boldsymbol{h})^{\mathsf{T}} (\boldsymbol{J}^2 + \lambda \boldsymbol{\Sigma})^{-1} (\frac{\tau_q}{\alpha}\boldsymbol{\Sigma}\boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}}\boldsymbol{\Sigma}^{1/2}\boldsymbol{h}) - \gamma^2 \right].$$

We compute the limit of the right hand side for any fixed value of $\lambda \ge 0$. First note that since $\|(\mathbf{\Sigma}^{-1/2} \mathbf{J}^2 \mathbf{\Sigma}^{-1/2} + \lambda \mathbf{I})^{-1}\| = O_p(1)$ and by invoking (E.29), the cross terms vanish in the limit and we have

$$\lim_{n\to\infty} \left(\frac{\tau_q}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \boldsymbol{h}\right)^{\mathsf{T}} \left(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^2 \boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I}\right)^{-1} \left(\frac{\tau_q}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \boldsymbol{h}\right)$$

(E.32)

$$= \lim_{n \to \infty} \frac{\tau_q^2}{\alpha^2} \boldsymbol{\theta}_0 \boldsymbol{\Sigma}^{1/2} (\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^2 \boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I})^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 + \lim_{n \to \infty} \frac{\beta^2}{d} \boldsymbol{h}^{\mathsf{T}} (\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^2 \boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I})^{-1} \boldsymbol{h}.$$

We treat each term separately. Plugging for θ_0 from (E.2) we have

$$\lim_{n \to \infty} \frac{\tau_q^2}{\alpha^2} \boldsymbol{\theta}_0 \boldsymbol{\Sigma}^{1/2} (\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^2 \boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I})^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 = \lim_{n \to \infty} (\frac{\mu_1 \tau_q}{\alpha})^2 \boldsymbol{\beta}^\mathsf{T} \boldsymbol{W}^\mathsf{T} (\boldsymbol{J}^2 + \lambda \boldsymbol{\Sigma})^{-1} \boldsymbol{W} \boldsymbol{\beta}$$
(E.33)
$$= \lim_{n \to \infty} (\frac{\mu_1 \tau_q}{\alpha})^2 \frac{1}{d} \operatorname{trace} (\boldsymbol{W}^\mathsf{T} (\boldsymbol{J}^2 + \lambda \boldsymbol{\Sigma})^{-1} \boldsymbol{W}),$$

where in the last step we used the fact that the distribution of W is rotationally invariant and $\|\beta\|_{\ell_0} \to 1$.

Similarly since $\boldsymbol{h} \sim N(0, \boldsymbol{I}_N)$ we have

$$\lim_{n \to \infty} \frac{\beta^2}{d} \boldsymbol{h}^{\mathsf{T}} (\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^2 \boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I})^{-1} \boldsymbol{h} = \lim_{n \to \infty} \frac{\beta^2}{d} \langle \boldsymbol{\Sigma}^{1/2} (\boldsymbol{J}^2 + \lambda \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma}^{1/2}, \boldsymbol{h} \boldsymbol{h}^{\mathsf{T}} \rangle$$

$$= \lim_{n \to \infty} \frac{\beta^2}{d} \operatorname{trace} (\boldsymbol{\Sigma}^{1/2} (\boldsymbol{J}^2 + \lambda \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma}^{1/2}).$$
(E.34)

Combining (E.33) and (E.34) into (E.32) we get

(E.35)
$$\lim_{n \to \infty} \left(\frac{\tau_q}{\alpha} \mathbf{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \boldsymbol{h} \right)^{\mathsf{T}} \left(\mathbf{\Sigma}^{-1/2} \boldsymbol{J}^2 \mathbf{\Sigma}^{-1/2} + \lambda \boldsymbol{I} \right)^{-1} \left(\frac{\tau_q}{\alpha} \mathbf{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \boldsymbol{h} \right)$$

$$= \lim_{n \to \infty} \operatorname{trace} \left\{ (\boldsymbol{J}^2 + \lambda \boldsymbol{\Sigma})^{-1} \left(\left(\frac{\mu_1 \tau_q}{\alpha} \right)^2 \frac{1}{d} \boldsymbol{W} \boldsymbol{W}^{\mathsf{T}} + \frac{\beta^2}{d} \boldsymbol{\Sigma} \right) \right\},$$

where we used that trace(AB) = trace(BA) for any two matrices A and B.

To calculate the limit on the right-hand side of (E.35), we use Proposition F.3 on the spectrum of inner product kernel random matrices.

Since $\|\mathbf{W}\| = O_p(1)$ we also have $\left\| \left(\frac{\mu_1 \tau_q}{\alpha} \right)^2 \frac{1}{d} \mathbf{W} \mathbf{W}^\mathsf{T} + \frac{\beta^2}{d} \mathbf{\Sigma} \right\| = O_p(1)$ and as an immediate corollary of Proposition F.3, in (E.35) we can replace \mathbf{J}^2 with \mathbf{K} since they have the same spectrum. This brings us to

$$\lim_{n \to \infty} \left(\frac{\tau_{q}}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} - \frac{\beta}{\sqrt{d}} \boldsymbol{h}\right)^{\mathsf{T}} \left(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^{2} \boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I}\right)^{-1} \left(\frac{\tau_{q}}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_{0} - \frac{\beta}{\sqrt{d}} \boldsymbol{h}\right)$$

$$= \lim_{n \to \infty} \operatorname{trace} \left\{ \left(\boldsymbol{K} + \lambda \boldsymbol{\Sigma}\right)^{-1} \left(\left(\frac{\mu_{1} \tau_{q}}{\alpha}\right)^{2} \cdot \frac{1}{d} \boldsymbol{W} \boldsymbol{W}^{\mathsf{T}} + \frac{\beta^{2}}{d} \boldsymbol{\Sigma}\right) \right\}$$
(E.36)
$$= \lim_{n \to \infty} \frac{1}{d} \operatorname{trace} \left\{ \left(\left(\frac{1}{4} + \lambda \mu_{1}^{2}\right) \boldsymbol{W} \boldsymbol{W}^{\mathsf{T}} + \left(\frac{1}{4} + \lambda \mu_{2}^{2}\right) \boldsymbol{I}\right)^{-1} \left(\left(\frac{\mu_{1}^{2} \tau_{q}^{2}}{\alpha^{2}} + \beta^{2} \mu_{1}^{2}\right) \boldsymbol{W} \boldsymbol{W}^{\mathsf{T}} + \beta^{2} \mu_{2}^{2} \boldsymbol{I}\right) \right\}.$$

Note that the latter only depends on the spectral density of WW^T and can be written in terms of its Stieltjes transform.

By law of large numbers and with simple algebraic manipulations it is easy to see that for any constants b_0, b_1, c_0, c_1 we have

(E.37)
$$\frac{1}{N} \sum_{i=1}^{N} \frac{b_0 s_i + b_1}{c_0 s_i + c_1} \to \frac{b_0}{c_0} + \frac{b_0 \frac{c_1}{c_0} - b_1}{c_0} S(-c_1/c_0; \psi_1),$$

with $S(t; \psi_1)$ representing the Stieltjes transform of the spectral density of WW^T . Using this with (E.36) we obtain

$$\lim_{n\to\infty} \left(\frac{\tau_q}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \boldsymbol{h}\right)^{\mathsf{T}} \left(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^2 \boldsymbol{\Sigma}^{-1/2} + \lambda \boldsymbol{I}\right)^{-1} \left(\frac{\tau_q}{\alpha} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \boldsymbol{h}\right)$$

(E.38)

$$=\frac{4\psi_1}{1+4\lambda\mu_1^2}\Big(\frac{\mu_1^2\tau_q^2}{\alpha^2}+\beta^2\mu_1^2\Big)+\frac{4\psi_1}{1+4\lambda\mu_1^2}\left\{\Big(\frac{\mu_1^2\tau_q^2}{\alpha^2}+\beta^2\mu_1^2\Big)\Big(\frac{1+4\lambda\mu_2^2}{1+4\lambda\mu_1^2}\Big)-\beta^2\mu_2^2\right\}S\Big(-\frac{1+4\lambda\mu_2^2}{1+4\lambda\mu_1^2};\psi_1\Big).$$

By combining (E.38) and (E.31) we get

$$\lim_{n\to\infty} Q(\boldsymbol{\Sigma}^{-1/2}\boldsymbol{J}, \frac{\tau_q}{\alpha}\boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{n}}\boldsymbol{h}, \gamma)$$

(E.39)

$$= \sup_{\lambda \geq 0} \quad \frac{\lambda}{2} \left[\frac{4\psi_1}{1 + 4\lambda\mu_1^2} \Big(\frac{\mu_1^2\tau_q^2}{\alpha^2} + \beta^2\mu_1^2 \Big) + \frac{4\psi_1}{1 + 4\lambda\mu_1^2} \left\{ \Big(\frac{\mu_1^2\tau_q^2}{\alpha^2} + \beta^2\mu_1^2 \Big) \Big(\frac{1 + 4\lambda\mu_2^2}{1 + 4\lambda\mu_1^2} \Big) - \beta^2\mu_2^2 \right\} S\Big(- \frac{1 + 4\lambda\mu_2^2}{1 + 4\lambda\mu_1^2}; \psi_1\Big) - \gamma^2 \right].$$

Plugging for $\mu_1 = \frac{1}{2}$ and $\mu_2 = \sqrt{\frac{1}{4} - \frac{1}{2\pi}}$ we have

(E.40)
$$\lim_{n \to \infty} Q(\mathbf{\Sigma}^{-1/2} \mathbf{J}, \frac{\tau_q}{\alpha} \mathbf{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{n}} \mathbf{h}, \gamma) = \mathsf{F}\left(\frac{\tau_q}{\alpha}, \beta, \psi_1, \gamma\right),$$

with the definition

(E.41)

$$\mathsf{F}(a,b,\psi_1,\gamma) \coloneqq \sup_{\lambda \geq 0} \ \frac{\lambda \psi_1}{2(1+\lambda)} \left\{ a^2 + b^2 + \left(a^2 \left(1 - \frac{2}{\pi} \frac{\lambda}{1+\lambda} \right) + \frac{2b^2}{\pi(1+\lambda)} \right) S\left(\frac{2}{\pi} \frac{\lambda}{1+\lambda} - 1; \psi_1 \right) \right\} - \frac{\lambda}{2} \gamma^2 \ .$$

By the change of variable $\tilde{\lambda} = \frac{\lambda}{1+\lambda}$, the function F can be written as

(E.42)

$$\mathsf{F}(a,b,\psi_1,\gamma)\coloneqq \sup_{0\leq \tilde{\lambda}<1} \ \frac{\tilde{\lambda}\psi_1}{2} \left\{ a^2 + b^2 + \left(a^2 \left(1 - \frac{2}{\pi} \tilde{\lambda} \right) + \frac{2(1-\tilde{\lambda})b^2}{\pi} \right) S\left(\frac{2}{\pi} \tilde{\lambda} - 1; \psi_1 \right) \right\} - \frac{\tilde{\lambda}}{2(1-\tilde{\lambda})} \gamma^2 \ .$$

We next proceed to characterize the limit of $\frac{1}{n}G_n(\boldsymbol{w} - \alpha \boldsymbol{g}; \frac{\tau_g}{\beta}, \gamma, \nu)$. To this end, we recall the result of [46, Lemma 6.4].

Lemma E.5 Let $u \in \mathbb{R}^n$ be a Gaussian random vector distributed as $N(\mathbf{0}, \omega^2 \mathbf{I}_n)$. Then,

(E.43)

$$\lim_{n\to\infty} \frac{1}{2n\rho(\rho+1)} \|\boldsymbol{u} - \mathsf{ST}(\boldsymbol{u};\nu)\|_{\ell_2}^2 = \frac{\omega^2}{2\rho(\rho+1)} \left(\left(1 - \sqrt{\frac{2}{\pi}} \frac{\nu}{\omega} e^{-\frac{\nu^2}{2\omega^2}}\right) \right)$$

(E.44)

$$\lim_{n\to\infty} \frac{1}{2n^2} \left(\frac{n}{\varepsilon} \gamma - \frac{1}{1+\rho} \left\| \mathsf{ST}(\boldsymbol{u}; \nu) \right\|_{\ell_1} \right)_+^2 = \frac{\omega^2}{2(\rho+1)^2} \left(\frac{\gamma(\rho+1)}{\varepsilon\omega} + \frac{\nu}{\omega} \cdot \operatorname{erfc}\left(\frac{1}{\sqrt{2}} \frac{\nu}{\omega} \right) - \sqrt{\frac{2}{\pi}} e^{-\frac{\nu^2}{2\omega^2}} \right)_+^2.$$

Therefore, by (E.22) we have

$$\lim_{n \to \infty} \frac{1}{n} G_n(\boldsymbol{u}; \rho, \gamma, \nu) = \frac{\omega^2}{2\rho(\rho + 1)} \left(\left(1 - \sqrt{\frac{2}{\pi}} \frac{\nu}{\omega} e^{-\frac{\nu^2}{2\omega^2}} \right) + \left(\frac{\nu^2}{\omega^2} - 1 \right) \operatorname{erfc} \left(\frac{1}{\sqrt{2}} \frac{\nu}{\omega} \right) \right) - \frac{\omega^2}{2(\rho + 1)^2} \left(\frac{\gamma(\rho + 1)}{\varepsilon \omega} + \frac{\nu}{\omega} \cdot \operatorname{erfc} \left(\frac{1}{\sqrt{2}} \frac{\nu}{\omega} \right) - \sqrt{\frac{2}{\pi}} e^{-\frac{\nu^2}{2\omega^2}} \right)_+^2.$$

Furthermore,

$$\min_{\nu \geq 0} \lim_{n \to \infty} \frac{1}{n} G_n(\boldsymbol{u}; \rho, \gamma, \nu)$$

$$=\mathsf{G}(\omega;\rho,\gamma) \coloneqq \begin{cases} 0 & \text{if } \gamma(\rho+1) \leq \sqrt{\frac{2}{\pi}}\varepsilon\omega \\ \frac{\omega^2}{2\rho(\rho+1)} \left(\operatorname{erf}\left(\frac{\nu^*\left(\frac{\gamma(\rho+1)}{\varepsilon\omega},\rho\right)}{\sqrt{2}}\right) - \frac{\gamma(\rho+1)}{\varepsilon\omega}\nu^*\left(\frac{\gamma(\rho+1)}{\varepsilon\omega},\rho\right) \right) & \text{if } \gamma(\rho+1) > \sqrt{\frac{2}{\pi}}\varepsilon\omega \end{cases}$$

where $\nu^*(a,\rho)$ is the unique solution to

$$a - \frac{1}{\rho}\nu - \nu \cdot \operatorname{erf}\left(\frac{\nu}{\sqrt{2}}\right) - \sqrt{\frac{2}{\pi}}e^{-\frac{\nu^2}{2}} = 0.$$

Combining (E.27), (E.29), (E.30), (E.40), and Lemma E.5, we obtain the following scalarized AO problem:

$$\max_{0 \le \beta \le K_{\beta}, 0 \le \gamma, \tau_{q}} \min_{0 \le \alpha < K_{\alpha}, 0 \le \tau_{g}} - \frac{\alpha}{\tau_{q}} \mathsf{F} \Big(\frac{\tau_{q}}{\alpha}, \beta, \psi_{1}, \gamma \Big) + \frac{\tau_{q}}{2\alpha} (\tau^{2} + 1 - \sigma^{2}) - \frac{\alpha \tau_{q}}{2} \\
+ \frac{\beta \tau_{g}}{2} \psi_{2} + \frac{\beta}{2(\tau_{q} + \beta)} (\sigma^{2} + \alpha^{2}) + \mathsf{G} \Big(\sqrt{\sigma^{2} + \alpha^{2}}; \frac{\tau_{g}}{\beta}, \gamma \Big),$$
(E.45)

where
$$\sigma^2 = \tau^2 + 1 - \psi_1 \left(1 + \left(1 - \frac{2}{\pi} \right) S\left(\frac{2}{\pi} - 1; \psi_1 \right) \right)$$
.

We conclude this part by a lemma on the convexity-concavity of the above scalarized AO problem and the uniqueness of the solution to the AO problem.

Lemma E.6 (Strict convexity and uniqueness of the solution) The objective function (E.45) is strictly jointly convex in (α, τ_g) and jointly concave in (β, γ, τ_q) . Also the solution $(\alpha_*, \frac{\tau_{g*}}{\beta_*})$ to this problem is unique.

We defer the proof of Lemma E.6 to Section E.5. This concludes the proof of Theorem 4.2(a).

E.2.2. Uniform convergence In Section E.2.1 we showed that the objective function in (E.24) converges point-wise to the objective function in (E.45). However, for our goal we need to show that the minimax solutions of the converging sequence of the objectives in (E.24) converges to the minimax solution of the AO objective in (E.45), denoted by $\mathcal{R}(\alpha, \tau_g, \beta, \gamma, \tau_q)$. Convexity/concavity of \mathcal{R} plays a crucial role here since it is being used to conclude local uniform convergence from the point-wise convergence.

This can be shown by following similar arguments as in [86, Lemma A.5] that is essentially based on a result known as "convexity lemma" in the literature (see e.g. [54, Lemma 7.75]) by which point-wise convergence of convex functions, of a finite number of variables, implies uniform convergence in compact subsets. Since the argument here is general, we leave out a detailed discussion and refer to [86, Lemma A.5].

E.3. Proof of Theorem 4.2(b) In Proposition 6.8 we gave a characterization of $\overset{\circ}{AR}_{nl}$. We first provide an alternative characterization in terms of the equivalent model of (E.2).

Recall the key quantity a from Proposition 6.8, given by

(E.46)
$$a^{2} = \tau^{2} + \left\| \frac{1}{2} \mathbf{W}^{\mathsf{T}} \boldsymbol{\theta} - \boldsymbol{\beta} \right\|_{\ell_{2}}^{2} + \left(\frac{1}{4} - \frac{1}{2\pi} \right) \|\boldsymbol{\theta}\|_{\ell_{2}}^{2}.$$

We claim that $a^2 = \sigma^2 + \|\mathbf{\Sigma}^{1/2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\|_{\ell_0}^2$. To see this, we expand this expression as follows:

$$\sigma^{2} + \|\boldsymbol{\Sigma}^{1/2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{0})\|_{\ell_{2}}^{2} = \sigma^{2} + \langle \boldsymbol{\theta}, \boldsymbol{\Sigma}\boldsymbol{\theta} \rangle + \langle \boldsymbol{\theta}_{0}, \boldsymbol{\Sigma}\boldsymbol{\theta}_{0} \rangle - 2\langle \boldsymbol{\theta}_{0}, \boldsymbol{\Sigma}\boldsymbol{\theta} \rangle$$

$$= \sigma^{2} + \mu_{1}^{2} \|\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\theta}\|_{\ell_{2}}^{2} + \mu_{2}^{2} \|\boldsymbol{\theta}\|_{\ell_{2}}^{2} + \mu_{1}^{2}\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}\boldsymbol{W}\boldsymbol{\beta} - 2\mu_{1}\langle \boldsymbol{W}, \boldsymbol{\beta}, \boldsymbol{\theta} \rangle$$

$$= \sigma^{2} + \|\mu_{1}\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\theta} - \boldsymbol{\beta}\|_{\ell_{2}}^{2} - \|\boldsymbol{\beta}\|_{\ell_{2}}^{2} + \mu_{2}^{2} \|\boldsymbol{\theta}\|_{\ell_{2}}^{2} + \mu_{1}^{2}\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}\boldsymbol{W}\boldsymbol{\beta}$$

$$= \tau^{2} + \|\mu_{1}\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\theta} - \boldsymbol{\beta}\|_{\ell_{2}}^{2} + \mu_{2}^{2} \|\boldsymbol{\theta}\|_{\ell_{2}}^{2},$$

where we used the definition of Σ , θ_0 and σ^2 as per (E.2). The claim follows by recalling that for the shifted Relu activation, $\mu_1 = \frac{1}{2}$ and $\mu_2 = \sqrt{\frac{1}{4} - \frac{1}{2\pi}}$. By the above characterization of quantity a we obtain

(E.47)
$$a^{2} = \sigma^{2} + \left\| \mathbf{\Sigma}^{1/2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{0}) \right\|_{\ell_{2}}^{2}.$$

We next note that by definition of the variables in the AO problem, we have $z = \Sigma^{1/2}(\theta - \theta_0)$ and $\alpha = \|\boldsymbol{z}\|_{\ell_2}$. Therefore,

$$\lim_{n\to\infty} \left\| \boldsymbol{\Sigma}^{1/2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \right\|_{\ell_2} = \alpha_* \,.$$

Invoking the limit of σ^2 given by (E.28), we get

(E.48)
$$\lim_{n \to \infty} a^2 = \tau^2 + 1 - \psi_1 \left(1 + \left(1 - \frac{2}{\pi} \right) S \left(\frac{2}{\pi} - 1 \right) \right) + \alpha_*^2.$$

We next characterize $\lim_{n\to\infty} \|J\theta\|_{\ell_2}$. We will use the same AO problem to calculate this

Recall that $\widehat{z} = \Sigma^{1/2} (\widehat{\theta}_{nl}^* - \theta_0)$ satisfies the following relation with q_* the optimizer in (E.13):

$$q_* = \arg\max_{\boldsymbol{q}} \boldsymbol{q}^\mathsf{T} \widehat{\boldsymbol{z}} - \widetilde{\ell}(\boldsymbol{v}; \boldsymbol{q}),$$

where $\tilde{\ell}(v;q)$ is the convex conjugate of $\bar{\ell}(v;z)$. Since conjugate of a conjugate function is the function itself we then have

$$\widehat{z} = \arg \max_{z} q_{*}^{\mathsf{T}} z - \overline{\ell}(v; z)$$

(E.49)

$$= \arg \max_{\boldsymbol{z}} \boldsymbol{q}_{\star}^{\mathsf{T}} \boldsymbol{z} - \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_{2}}^{2} - \frac{\varepsilon}{n} \|\boldsymbol{v}\|_{\ell_{1}} \|\boldsymbol{J}(\boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z})\|_{\ell_{2}} - \frac{\varepsilon^{2}}{2} \|\boldsymbol{J}(\boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \boldsymbol{z})\|_{\ell_{2}}^{2}.$$

We consider two cases:

Case 1: $\|J(\theta_0 + \Sigma^{-1/2}\widehat{z})\|_{\ell_2} \neq 0$. Setting derivative with respect to \widehat{z} to zero we obtain

$$\boldsymbol{q}_{*} - \frac{\varepsilon}{n} \|\boldsymbol{v}\|_{\ell_{1}} \frac{\boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^{2} (\boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \widehat{\boldsymbol{z}})}{\|\boldsymbol{J} (\boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \widehat{\boldsymbol{z}})\|_{\ell_{2}}} - \varepsilon^{2} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}^{2} (\boldsymbol{\theta}_{0} + \boldsymbol{\Sigma}^{-1/2} \widehat{\boldsymbol{z}}) = 0.$$

By rearranging the terms we write it as

$$J(\boldsymbol{\theta}_0 + \boldsymbol{\Sigma}^{-1/2}\widehat{\boldsymbol{z}}) = \left(\frac{\varepsilon}{n} \frac{\|\boldsymbol{v}\|_{\ell_1}}{\|J(\boldsymbol{\theta}_0 + \boldsymbol{\Sigma}^{-1/2}\widehat{\boldsymbol{z}})\|_{\ell_2}} + \varepsilon^2\right)^{-1} J^{-1}\boldsymbol{\Sigma}^{1/2}\boldsymbol{q}_*$$

By taking the ℓ_2 norm of both sides and then solving for $\|J(\theta_0 + \Sigma^{-1/2}\widehat{z})\|_{\ell_2}$, we get

(E.50)
$$\| \boldsymbol{J}(\boldsymbol{\theta}_0 + \boldsymbol{\Sigma}^{-1/2} \widehat{\boldsymbol{z}}) \|_{\ell_2} = \frac{1}{\varepsilon^2} \| \boldsymbol{J}^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{q}_* \|_{\ell_2} - \frac{1}{n\varepsilon} \| \boldsymbol{v} \|_{\ell_1} .$$

Case 2: $\|J(\theta_0 + \Sigma^{-1/2}\widehat{z})\|_{\ell_2} = 0$. In this case, $\widehat{z} = -\Sigma^{1/2}\theta_0$ and by comparing the objective function of (E.49) at the optimal solution under case 1 and case 2, it is easy to verify that case 2 happens only when the right-hand side in (E.50) becomes negative. Therefore, the two cases can be combined together in the following form:

(E.51)
$$\langle \widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^*, \boldsymbol{J}^2 \widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^* \rangle = \| \boldsymbol{J} (\boldsymbol{\theta}_0 + \boldsymbol{\Sigma}^{-1/2} \widehat{\boldsymbol{z}}) \|_{\ell_2}^2 = \left(\frac{1}{\varepsilon^2} \| \boldsymbol{J}^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{q}_* \|_{\ell_2} - \frac{1}{n\varepsilon} \| \boldsymbol{v} \|_{\ell_1} \right)_{\perp}^2.$$

So in order to get the asymptotic value of the left hand side we can work with the right-hand side with v and $\gamma = \| \boldsymbol{J}^{-1} \boldsymbol{\Sigma}^{1/2} \boldsymbol{q}_* \|_{\ell_2}$ the optimal solutions of the AO problem.

In Lemma E.4 (which is a restatement of [46, Lemma 6.3]), the Moreau envelop function $e_f(x, \rho)$ was characterized. Following the proof of [46, Lemma 6.3], we can verify that the optimal v is given by

$$v = \frac{1}{1 + \frac{\tau_g}{\beta}} ST(w - \alpha g; \nu).$$

Therefore, by invoking the relation (E.44) we have

$$\lim_{n\to\infty} \frac{1}{n^2} \left(\frac{n}{\varepsilon} \gamma - \|\boldsymbol{v}\|_{\ell_1} \right)_+^2 = \frac{\omega^2}{(\rho+1)^2} \left(\frac{\gamma(\rho+1)}{\varepsilon\omega} + \nu_* \cdot \operatorname{erfc} \left(\frac{1}{\sqrt{2}} \nu_* \right) - \sqrt{\frac{2}{\pi}} e^{-\nu_*^2} \right)_+^2,$$

with $\omega = \sqrt{\alpha^2 + \sigma^2}$, $\rho = \frac{\tau_g}{\beta}$, $\nu_* = \nu_*(\frac{\gamma(\rho+1)}{\varepsilon\omega}, \rho)$ and $\nu_*(a, \mu)$ the unique solution to the following equation:

$$a - \frac{1}{\rho}\nu - \nu \cdot \operatorname{erf}\left(\frac{\nu}{\sqrt{2}}\right) - \sqrt{\frac{2}{\pi}}e^{-\frac{\nu^2}{2}} = 0.$$

Plugging the above relation into (E.51) we obtain

$$\begin{split} \lim_{n \to \infty} \langle \widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^*, \boldsymbol{J}^2 \widehat{\boldsymbol{\theta}}_{\mathrm{nl}}^* \rangle &= \lim_{n \to \infty} \left(\frac{1}{\varepsilon^2} \gamma - \frac{1}{n\varepsilon} \| \boldsymbol{v} \|_{\ell_1} \right)_+^2 \\ &= \frac{1}{\varepsilon^2 n^2} \lim_{n \to \infty} \left(\frac{n}{\varepsilon} \gamma - \| \boldsymbol{v} \|_{\ell_1} \right)_+^2 \\ &= \frac{\omega^2}{\varepsilon^2 (\rho + 1)^2} \left(\frac{\gamma (\rho + 1)}{\varepsilon \omega} + \nu_* \cdot \operatorname{erfc} \left(\frac{1}{\sqrt{2}} \nu_* \right) - \sqrt{\frac{2}{\pi}} e^{-\nu_*^2} \right)_+^2 \\ &= \frac{\omega^2}{\varepsilon^2 (\rho + 1)^2} \left(\frac{\gamma (\rho + 1)}{\varepsilon \omega} + \frac{\rho + 1}{\rho} \nu_* - \frac{\gamma (\rho + 1)}{\varepsilon \omega} \right)_+^2 \\ &= \frac{\omega^2}{\varepsilon^2 (\rho + 1)^2} \left(\frac{1 + \rho}{\rho} \nu_* \right)_+^2 \end{split}$$

$$= \frac{\omega^2 \nu_*^2}{\varepsilon^2 \rho^2}$$

$$= \frac{(\alpha_*^2 + \sigma^2) \nu_*^2}{\varepsilon^2 \rho^2}$$

$$= \frac{\beta^2 \nu_*^2 (\alpha_*^2 + \sigma^2)}{\varepsilon^2 \tau_q^2}.$$
(E.52)

Now by recalling the characterization (6.20) along with (E.52) and (E.47) we get the desired result of (4.6).

E.4. Proof of Proposition 5.1 Recall the objective $\mathcal{R}(\alpha,\tau_g,\beta,\gamma,\tau_q)$. We start by considering the change of variable $\tilde{\gamma}=\gamma/\varepsilon$. Note that the term $-\lambda/(1-\lambda)\gamma^2=-\lambda/(1-\lambda)\varepsilon^2\tilde{\gamma}^2$ will be dropped as it is zero. We next argue that $\nu^*=0$. The reason is that if the indicator in the objective is inactive then the corresponding term is void, which is equivalent to $\nu^*=0$. If the indicator is active, since the problem is maximization over γ , we deduce that $\frac{\gamma(\tau_g+\beta)}{\varepsilon\beta\sqrt{\alpha^2+\sigma^2}}=\sqrt{\frac{2}{\pi}}$, which by the equation defining ν^* implies that $\nu^*=0$. Next, by straightforward calculation, and using definition of Stieltjes transform S, we have that the expression inside $\sup_{0\leq \lambda<1}$ is increasing in λ and so we have that the optimal $\lambda\to 1$. Using these values, the objective reduces to

$$\mathcal{R}(\alpha, \tau_g, \beta, \tau_q) = \frac{\tau_q}{2\alpha} (\tau^2 + 1 - \sigma^2) - \frac{\alpha \tau_q}{2} + \frac{\beta \tau_g}{2} \psi_2 + \frac{\beta}{2(\tau_g + \beta)} (\sigma^2 + \alpha^2)$$
$$- \frac{\psi_1}{2} \left\{ \frac{\tau_q}{\alpha} + \frac{\alpha}{\tau_q} \beta^2 + \frac{\tau_q}{\alpha} \left(1 - \frac{2}{\pi} \right) S\left(\frac{2}{\pi} - 1; \psi_1 \right) \right\}.$$

Using the definition

$$\sigma^2 = \tau^2 + 1 - \psi_1 \left(1 + \left(1 - \frac{2}{\pi} \right) S\left(\frac{2}{\pi} - 1; \psi_1 \right) \right),$$

we can further simplify the objective as

$$\mathcal{R}(\alpha, \tau_g, \beta, \tau_q) = -\frac{\psi_1}{2} \frac{\beta^2 \alpha}{\tau_q} - \frac{\alpha \tau_q}{2} + \frac{\beta \tau_g}{2} \psi_2 + \frac{\beta}{2(\tau_g + \beta)} (\sigma^2 + \alpha^2).$$

Optimization over τ_q can be done easily resulting in $\tau_q = \beta \sqrt{\psi_1}$, which gives

$$\mathcal{R}(\alpha, \tau_g, \beta) = -\alpha\beta\sqrt{\psi_1} + \frac{\beta\tau_g}{2}\psi_2 + \frac{\beta}{2(\tau_g + \beta)}(\sigma^2 + \alpha^2).$$

Writing the stationary condition for α, τ_q, β we arrive at the following system of equations:

(E.53)
$$\begin{cases} \frac{\alpha}{\tau_g + \beta} = \sqrt{\psi_1}, \\ \psi_2 = \frac{\sigma^2 + \alpha^2}{(\tau_g + \beta)^2}, \\ -\alpha\sqrt{\psi_1} + \frac{\tau_g\psi_2}{2} + \frac{\tau_g}{2} \frac{\sigma^2 + \alpha^2}{(\tau_g + \beta)^2} = 0. \end{cases}$$

Solving the above system of equations we obtain $\alpha^2 = \sigma^2 \psi_1 / (\psi_2 - \psi_1)$. Recalling that ν^* , using Theorem 4.2 (b) we get the standard risk of the estimator to be

$$SR(\widehat{\boldsymbol{\theta}}) = \alpha_*^2 + \sigma^2 = \sigma^2 \left(\frac{\psi_2}{\psi_2 - \psi_1} \right).$$

E.5. Proofs of the Auxiliary Lemmas

E.5.1. *Proof of Lemma E.1* We start by considering the following related but different function

$$\ell_0(\boldsymbol{v};\boldsymbol{z}) = \frac{1}{2n} \sum_{i=1}^n \left(|v_i| + \varepsilon \|\boldsymbol{z}\|_{\ell_2} \right)^2.$$

As shown in the proof of Lemma 6.1 in [46], the conjugate of this function is given by

$$\ell_0^*(\boldsymbol{v}; \boldsymbol{q}) = \frac{1}{2} \left(\frac{\|\boldsymbol{q}\|_{\ell_2}}{\varepsilon} - \frac{\|\boldsymbol{v}\|_{\ell_1}}{n} \right)_+^2 - \frac{1}{2n} \|\boldsymbol{v}\|_{\ell_2}^2.$$

Note that $\bar{\ell}(v; z) = \ell_0(v; J(\theta_0 + \Sigma^{-1/2}z))$. We next use the result that if $f(x) = g(Ax + x_0)$ then the conjugate of f can be written in terms of the conjugate of g as follows:

$$f^*(\boldsymbol{y}) = -\langle \boldsymbol{A}^{-1} \boldsymbol{x}_0, \boldsymbol{y} \rangle + g^*(\boldsymbol{A}^{-\mathsf{T}} \boldsymbol{y}).$$

Using this result with $x_0 = J\theta_0$ and $A = J\Sigma^{-1/2}$ we obtain

$$\widetilde{\ell}(\boldsymbol{v};\boldsymbol{q}) = -\langle \boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_0, \boldsymbol{q} \rangle + \frac{1}{2} \left(\frac{1}{\varepsilon} \left\| \boldsymbol{J}^{-1}\boldsymbol{\Sigma}^{1/2}\boldsymbol{q} \right\|_{\ell_2} - \frac{\|\boldsymbol{v}\|_{\ell_1}}{n} \right)_+^2 - \frac{1}{2n} \left\| \boldsymbol{v} \right\|_{\ell_2}^2$$

E.5.2. *Proof of Lemma E.2* Consider a slightly different optimization than (E.15) where the equality constraint is replaced by the inequality constraint $\|\widetilde{q}\|_{\ell_2} \le \gamma$:

(E.54)
$$\min_{\widetilde{\boldsymbol{q}},0 \le \gamma} \frac{c_0}{2} \|\boldsymbol{H}\widetilde{\boldsymbol{q}} - \boldsymbol{r}\|_{\ell_2}^2 + \frac{1}{2} \left(\frac{1}{\varepsilon}\gamma - c_1\right)_+^2$$
s.t. $\|\widetilde{\boldsymbol{q}}\|_{\ell_2} \le \gamma$.

Due to this change, (E.54) is now a convex optimization. Denote by OPT_1 the optimal objective value of the original problem (E.15) and by OPT_2 the optimal objective value of the modified problem (E.54). We argue that $\mathsf{OPT}_1 = \mathsf{OPT}_2$. Clearly $\mathsf{OPT}_1 \ge \mathsf{OPT}_2$ because (E.54) has a larger feasible set. Now suppose that this inequality is strict ($\mathsf{OPT}_1 > \mathsf{OPT}_2$) and let $(\widetilde{q}_*, \gamma_*)$ be a solution to (E.54). Then we should have $\|\widetilde{q}_*\|_{\ell_2} < \gamma_*$. Consider the point $(\widetilde{q}_*, \|\widetilde{q}_*\|_{\ell_2})$ which is a feasible point for both optimization problems and so the objective value at this point is at least OPT_1 and therefore strictly larger than OPT_2 . But this is a contradiction because $\left(\frac{1}{\varepsilon}\gamma - c_1\right)^2$ is non-decreasing in $\gamma \ge 0$.

To characterize OPT₂, we first focus on the minimization over \tilde{q} . The corresponding Lagrangian with Lagrange multiplier $\frac{\lambda c_0}{2}$ reads

$$\sup_{\lambda>0} \min_{\widetilde{\boldsymbol{q}}} \frac{c_0}{2} \|\boldsymbol{H}\widetilde{\boldsymbol{q}} - \boldsymbol{r}\|_{\ell_2}^2 - \frac{\lambda c_0}{2} (\gamma^2 - \|\widetilde{\boldsymbol{q}}\|_{\ell_2}^2).$$

Solving the inner minimization, we have $\widetilde{q}_* = (H^T H + \lambda I)^{-1} H^T r$ and the dual problem becomes

$$\sup_{\lambda \geq 0} \frac{c_0}{2} \|\boldsymbol{H}\widetilde{\boldsymbol{q}}_* - \boldsymbol{r}\|_{\ell_2}^2 - \frac{\lambda c_0}{2} (\gamma^2 - \|\widetilde{\boldsymbol{q}}_*\|_{\ell_2}^2)$$

$$= \sup_{\lambda \geq 0} \frac{c_0}{2} \widetilde{\boldsymbol{q}}_*^{\mathsf{T}} [\boldsymbol{H}^{\mathsf{T}} (\boldsymbol{H}\widetilde{\boldsymbol{q}}_* - \boldsymbol{r}) + \lambda \widetilde{\boldsymbol{q}}_*] - \frac{c_0}{2} \boldsymbol{r}^{\mathsf{T}} (\boldsymbol{H}\widetilde{\boldsymbol{q}}_* - \boldsymbol{r}) - \frac{\lambda c_0}{2} \gamma^2$$

$$= \sup_{\lambda \geq 0} -\frac{c_0}{2} \boldsymbol{r}^{\mathsf{T}} (\boldsymbol{H}\widetilde{\boldsymbol{q}}_* - \boldsymbol{r}) - \frac{\lambda c_0}{2} \gamma^2$$

$$= \sup_{\lambda \ge 0} -\frac{c_0}{2} \boldsymbol{r}^{\mathsf{T}} (\boldsymbol{H} (\boldsymbol{H}^{\mathsf{T}} \boldsymbol{H} + \lambda \boldsymbol{I})^{-1} \boldsymbol{H}^{\mathsf{T}} - \boldsymbol{I}) \boldsymbol{r} - \frac{\lambda c_0}{2} \gamma^2$$

$$= \sup_{\lambda \ge 0} \frac{\lambda c_0}{2} \boldsymbol{r}^{\mathsf{T}} (c_0 \boldsymbol{H}^{\mathsf{T}} \boldsymbol{H} + \lambda \boldsymbol{I})^{-1} \boldsymbol{r} - \frac{\lambda c_0}{2} \gamma^2$$
(E.55)
$$= c_0 Q(\boldsymbol{H}, \boldsymbol{r}, \gamma).$$

By the Slater's condition the duality gap is zero and hence by next minimizing over $\gamma \ge 0$, we obtain that the optimal value of (E.54) is given by

$$\min_{\gamma \geq 0} c_0 Q(\boldsymbol{H}, \boldsymbol{r}, \gamma) + \frac{1}{2} \left(\frac{1}{\varepsilon} \gamma - c_1 \right)_+^2.$$

E.5.3. *Proof of Lemma E.3* We first show that the function

$$g(\gamma, \beta) = Q(\mathbf{\Sigma}^{-1/2} \mathbf{J}, \frac{1}{\alpha} \mathbf{\Sigma}^{1/2} \boldsymbol{\theta}_0 - \frac{\beta}{\sqrt{d}} \mathbf{h}, \gamma)$$

is jointly convex in (γ, β) . By the change of variable $\tilde{\lambda} = \lambda \gamma$, $\tilde{\boldsymbol{\theta}} = \boldsymbol{\Sigma}^{1/2} \boldsymbol{\theta}_0 / \alpha$, $\boldsymbol{H} = \boldsymbol{\Sigma}^{-1/2} \boldsymbol{J}$, this function can be written as

$$g(\gamma,\beta) = \sup_{\tilde{\lambda} \geq 0} \frac{\tilde{\lambda}}{2\gamma} \left((\tilde{\boldsymbol{\theta}} - \frac{\beta}{\sqrt{d}} \boldsymbol{h})^{\mathsf{T}} (\boldsymbol{H} \boldsymbol{H}^{\mathsf{T}} + \frac{\tilde{\lambda}}{\gamma} \boldsymbol{I})^{-1} (\tilde{\boldsymbol{\theta}} - \frac{\beta}{\sqrt{d}} \boldsymbol{h}) - \gamma^{2} \right)$$
$$= \sup_{\tilde{\lambda} > 0} \frac{\tilde{\lambda}}{2} \left((\tilde{\boldsymbol{\theta}} - \frac{\beta}{\sqrt{d}} \boldsymbol{h})^{\mathsf{T}} (\gamma \boldsymbol{H} \boldsymbol{H}^{\mathsf{T}} + \tilde{\lambda} \boldsymbol{I})^{-1} (\tilde{\boldsymbol{\theta}} - \frac{\beta}{\sqrt{d}} \boldsymbol{h}) - \gamma \right).$$

We show that for any fixed $\tilde{\lambda} \ge 0$ the inner function above is jointly convex in (γ, β) and since the pointwise maximum of convex functions is also convex, we conclude that $g(\gamma, \beta)$ is jointly convex in (γ, β) .

The Hessian of the inner function reads

$$\frac{1}{2}\nabla^{2}_{\frac{\beta}{\sqrt{d}},\gamma}\left[(\tilde{\boldsymbol{\theta}}-\frac{\beta}{\sqrt{d}}\boldsymbol{h})^{\mathsf{T}}(\gamma\boldsymbol{H}\boldsymbol{H}^{\mathsf{T}}+\tilde{\lambda}\boldsymbol{I})^{-1}(\tilde{\boldsymbol{\theta}}-\frac{\beta}{\sqrt{d}}\boldsymbol{h})-\gamma\right]=\begin{bmatrix}A & C\\ C & B\end{bmatrix},$$

where

$$A := \mathbf{h}^{\mathsf{T}} (\gamma \mathbf{H} \mathbf{H}^{\mathsf{T}} + \tilde{\lambda} \mathbf{I})^{-1} \mathbf{h}$$

$$B := \left\| (\gamma \mathbf{H} \mathbf{H}^{\mathsf{T}} + \tilde{\lambda} \mathbf{I})^{-1/2} \mathbf{H} \mathbf{H}^{\mathsf{T}} (\gamma \mathbf{H} \mathbf{H}^{\mathsf{T}} + \tilde{\lambda} \mathbf{I})^{-1} (\tilde{\boldsymbol{\theta}} - \frac{\beta}{\sqrt{d}} \mathbf{h}) \right\|_{\ell_{2}}^{2}$$

$$C := \mathbf{h}^{\mathsf{T}} (\gamma \mathbf{H} \mathbf{H}^{\mathsf{T}} + \tilde{\lambda} \mathbf{I})^{-1} \mathbf{H} \mathbf{H}^{\mathsf{T}} (\gamma \mathbf{H} \mathbf{H}^{\mathsf{T}} + \tilde{\lambda} \mathbf{I})^{-1} (\tilde{\boldsymbol{\theta}} - \frac{\beta}{\sqrt{d}} \mathbf{h}).$$

Here we repeatedly used the identity $\frac{\partial \mathbf{K}^{-1}}{\partial \gamma} = -\mathbf{K}^{-1} \frac{\partial \mathbf{K}}{\partial \gamma} \mathbf{K}^{-1}$, for a matrix \mathbf{K} .

To lighten the notation, we set $M := (\gamma H H^{\mathsf{T}} + \tilde{\lambda} I)^{-1}$ and $v := \tilde{\theta} - \frac{\beta}{\sqrt{d}} h$. Using these shorthands the determinant of the Hessian is equal to

$$\left\|\boldsymbol{M}^{1/2}\boldsymbol{h}\right\|_{\ell_{2}}^{2}\left\|\boldsymbol{M}^{1/2}\boldsymbol{H}\boldsymbol{H}^{\mathsf{T}}\boldsymbol{M}\boldsymbol{v}\right\|_{\ell_{2}}^{2}-(\boldsymbol{h}^{\mathsf{T}}\boldsymbol{M}\boldsymbol{H}\boldsymbol{H}^{\mathsf{T}}\boldsymbol{M}\boldsymbol{v})^{2}\geq0,$$

using the Cauchy–Schwarz inequality. This completes the proof of $g(\gamma, \beta)$ being jointly convex in (γ, β) .

Next note that its perspective function is given by

$$\tau_{q}g(\gamma/\tau_{q},\beta/\tau_{q}) = \tau_{q}Q(\boldsymbol{\Sigma}^{-1/2}\boldsymbol{J}, \frac{1}{\alpha}\boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_{0} - \frac{\beta}{\tau_{q}\sqrt{d}}\boldsymbol{h}, \frac{\gamma}{\tau_{q}})$$

$$= \frac{1}{\tau_{q}}Q(\boldsymbol{\Sigma}^{-1/2}\boldsymbol{J}, \frac{\tau_{q}}{\alpha}\boldsymbol{\Sigma}^{1/2}\boldsymbol{\theta}_{0} - \frac{\beta}{\sqrt{d}}\boldsymbol{h}, \gamma),$$

and therefore is jointly convex in (γ, β, τ_q) .

E.5.4. Proof of Lemma E.6 As we discussed after (E.24), the objective function in (E.24) is jointly convex in (α, τ_g) and jointly concave in (β, γ, τ_q) . Since convexity/concavity is preserved by point-wise limits, the objective (E.45) is jointly convex in (α, τ_g) and jointly concave in (β, γ, τ_q) . To prove strict convexity in (α, τ_g) , note that in our derivation we wrote (E.19) (the part of the objective E.18 that involves v) in terms of the Moreau envelope $\frac{1}{n}e_f\left(w-\alpha g;\frac{\tau_g}{\beta}\right)$, cf. (E.21). As $d\to\infty$, its limit goes to the expected Moreau envelope. By using the result of [86, Lemma 4.4] the expected Moreau envelope of a function is strictly convex in $\mathbb{R}_{>0}\times\mathbb{R}_{>0}$ without requiring any strong or strict convexity assumption on the function itself. Therefore, the objective (E.24) (and so objective of (E.45) after taking point-wise limit) is jointly strictly convex in (α, τ_g) .

To prove the uniqueness, note that $\max_{0 \le \beta, \gamma, \tau_q} \mathcal{R}(\alpha, \tau_g, \beta, \gamma, \tau_q)$ is strictly convex in (α, τ_g) . This follows from the fact that if $f(\boldsymbol{x}, \boldsymbol{y})$ is strictly convex in \boldsymbol{x} , then $\max_{\boldsymbol{y}} f(\boldsymbol{x}, \boldsymbol{y})$ is also strictly convex in \boldsymbol{x} . We next use [86, Lemma C.5] to conclude that $\min_{\tau_g > 0} \max_{0 \le \beta, \gamma, \tau_q} \mathcal{R}(\alpha, \tau_g, \beta, \gamma, \tau_q)$ is strictly convex in $\alpha \ge 0$. Therefore, its minimizer over $\alpha \ge 0$ is unique. By a similar argument, we show that $\frac{\tau_{g*}}{\beta_*}$ is unique. Consider the change of variable $\tau_g \to \tilde{\tau}_g = \frac{\tau_g}{\beta}$. Then part of the objective (E.24) that depends on $\tilde{\tau}_g$ can be written as

$$\frac{\beta^{2}\tilde{\tau_{g}}}{2}\frac{n}{d} + \frac{1}{2(\tilde{\tau_{g}}+1)}\frac{1}{n} \|\boldsymbol{w} - \alpha \boldsymbol{g}\|_{\ell_{2}}^{2} + \frac{1}{n}G_{n}(\boldsymbol{w} - \alpha \boldsymbol{g}; \tilde{\tau_{g}}, \gamma, \nu) = \frac{\beta^{2}\tilde{\tau_{g}}}{2}\frac{n}{d} + \frac{1}{n}e_{f}(\boldsymbol{w} - \alpha \boldsymbol{g}; \tilde{\tau_{g}}),$$

using (E.19). As explained above, this converges to the expected Moreau envelope, which is strictly convex in $\tilde{\tau_g}$. Following by the same reasoning for α , one can show that $\min_{\alpha>0} \max_{0\leq \beta,\gamma,\tau_q} \mathcal{R}(\alpha,\tilde{\tau_g},\beta,\gamma,\tau_q)$ is strictly convex in $\tilde{\tau_g}>0$. Therefore, its minimizer over $\tilde{\tau_g}>0$ is unique.

APPENDIX F: SOME USEFUL LEMMAS

Here we state some of the technical lemmas that are used in deriving our analytical results. The first lemma is about the Stieltjes transform of the Marchenko-Pastur distribution.

Definition F.1 The Stieltjes transform $S_{\rho}(z)$ of a measure of density ρ on a real interval I is the function of the complex variable z defined outside I by the formula

$$S_{\rho}(z) = \int_{I} \frac{\rho(t) dt}{z - t} \quad z \in \mathbb{C} \backslash I.$$

Lemma F.2 Suppose that $\mathbf{W} \in \mathbb{R}^{N \times d}$ has rows drawn independently from unit sphere. As $N, d \to \infty$ and $N/d \to \psi_1$, the spectral density of $\mathbf{W}\mathbf{W}^\mathsf{T}$ converges (in weak topology in distribution) to the Marchenko-Pastur distribution with Stieltjes transform given by

$$S(z;\psi_1) = \frac{1 - \psi_1 - z - \sqrt{(1 - \psi_1 - z)^2 - 4\psi_1 z}}{-2\psi_1 z},$$

Proof We refer to [2, page 52] for the proof of this proposition.

The next lemma is about the spectrum of matrix J given by

(F.1)
$$J = \left(W W^{\mathsf{T}} \right) \odot \left(\frac{\pi - \cos^{-1} (W W^{\mathsf{T}})}{2\pi} \right)^{1/2}.$$

Lemma F.3 Suppose that $W \in \mathbb{R}^{N \times d}$ has rows chosen randomly and independently of data form the unit sphere, $\mathrm{Unif}(\mathbb{S}^{d-1})$. Let J be given by (F.1) and suppose that $N/d \to \psi_1 \in (0, \infty)$, as $n \to \infty$. Then, the matrix J^2 can (in probability) be approximated consistently in operator norm by the matrix K given by

$$\boldsymbol{K} = \frac{1}{4}(\boldsymbol{W}\boldsymbol{W}^{\mathsf{T}} + \boldsymbol{I}).$$

In other words, $\|\mathbf{J}^2 - \mathbf{K}\| \to 0$, in probability, when $n \to \infty$.

Proof The claim follows from the result of [26, Theorem 2.1] about the spectrum of inner product kernel random matrices, specialized to matrix J^2 . Specifically, let $f(z) = z(\pi - \cos^{-1}(z))/(2\pi)$. Then $J_{ij}^2 = f(w_i^{\mathsf{T}}w_j)$. By employing [26, Theorem 2.1], the kernel matrix J^2 can (in probability) be approximated consistently in operator norm by the matrix K, given by

$$K = f(0)\mathbf{1}\mathbf{1}^{\mathsf{T}} + f'(0)WW^{\mathsf{T}} + (f(1) - f(0) - f'(0))I.$$

For our specific f we have f(0) = 0, f(1) = 1/2, f'(0) = 1/4.