

Article

# No Statistical-Computational Gap in Spiked Matrix Models with Generative Network Priors<sup>†</sup>

Jorio Cocola<sup>1,\*</sup>, Paul Hand<sup>1,2</sup> and Vladislav Voroninski<sup>3</sup>

<sup>1</sup> Department of Mathematics, Northeastern University, Boston, MA 02115, USA

<sup>2</sup> Khoury College of Computer Sciences, Northeastern University, Boston, MA 02115, USA; p.hand@northeastern.edu

<sup>3</sup> Helm.ai, Menlo Park, CA 94025, USA; vlad@helm.ai

\* Correspondence: cocola.j@northeastern.edu

† This paper is an extended version of our paper published in NeurIPS 2020.

**Abstract:** We provide a non-asymptotic analysis of the spiked Wishart and Wigner matrix models with a generative neural network prior. Spiked random matrices have the form of a rank-one signal plus noise and have been used as models for high dimensional Principal Component Analysis (PCA), community detection and synchronization over groups. Depending on the prior imposed on the spike, these models can display a statistical-computational gap between the information theoretically optimal reconstruction error that can be achieved with unbounded computational resources and the sub-optimal performances of currently known polynomial time algorithms. These gaps are believed to be fundamental, as in the emblematic case of Sparse PCA. In stark contrast to such cases, we show that there is no statistical-computational gap under a generative network prior, in which the spike lies on the range of a generative neural network. Specifically, we analyze a gradient descent method for minimizing a nonlinear least squares objective over the range of an expansive-Gaussian neural network and show that it can recover in polynomial time an estimate of the underlying spike with a rate-optimal sample complexity and dependence on the noise level.

**Keywords:** spiked matrix models; generative networks; rank-one matrix recovery; statistical-computational gap



**Citation:** Cocola, J.; Hand, P.; Voroninski, V. No Statistical-Computational Gap in Spiked Matrix Models with Generative Network Priors. *Entropy* **2021**, *23*, 115. <https://doi.org/10.3390/e23010115>

Received: 1 December 2020

Accepted: 8 January 2021

Published: 16 January 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

One of the fundamental problems in statistical inference and signal processing is the estimation of a signal given noisy high dimensional data. A prototypical example is provided by spiked matrix models where a signal  $y_* \in \mathbb{R}^n$  is to be estimated from a matrix  $Y$  taking one of the following forms:

- **Spiked Wishart Model** in which  $Y \in \mathbb{R}^{N \times n}$  is given by:

$$Y = u y_*^T + \sigma Z, \quad (1)$$

where  $\sigma > 0$ ,  $u \sim \mathcal{N}(0, I_N)$ ,  $Z_{ij}$  are i.i.d. from  $\mathcal{N}(0, 1)$ , and  $u$  and  $Z$  are independent;

- **Spiked Wigner Model** in which  $Y \in \mathbb{R}^{n \times n}$  is given by:

$$Y = y_* y_*^T + \nu \mathcal{H} \quad (2)$$

where  $\nu > 0$ ,  $\mathcal{H} \in \mathbb{R}^{n \times n}$  is drawn from a *Gaussian Orthogonal Ensemble*  $\text{GOE}(n)$ , that is,  $\mathcal{H}_{ii} \sim \mathcal{N}(0, 2/n)$  for all  $1 \leq i \leq n$  and  $\mathcal{H}_{ij} = \mathcal{H}_{ji} \sim \mathcal{N}(0, 1/n)$  for  $1 \leq j < i \leq n$ .

In the last 20 years, spiked random matrices have been extensively studied, as they serve as a mathematical model for many signal recovery problems such as PCA [1–4], synchronization over graphs [5–7] and community detection [8–10]. Furthermore, these models are archetypal examples of the trade-off between statistical accuracy and computational efficiency. From a statistical perspective, the objective is to understand how the

choice of the prior on  $y_*$  determines the critical signal-to-noise ratio (SNR) and number of measurements above which it becomes information-theoretically possible to estimate the signal. From a computational perspective, the objective is to design efficient algorithms that leverage such prior information. A recent and vast body of literature has shown that depending on the chosen prior, gaps can arise between the minimum SNR required to solve the problem and the one above which known polynomial-time algorithms succeed. An emblematic example is provided the Sparse PCA problem where the signal  $y_*$  in (1) is taken to be  $s$ -sparse. In this case  $N = \mathcal{O}(s \log n)$  number of samples are sufficient for estimating  $y_*$  [2,4], while the best known efficient algorithms require  $N = \mathcal{O}(s^2)$  [3,11,12]. This gap is believed to be fundamental. This “statistical-computational gap” has been observed also for Spiked Wigner models (2) and, in general, for other structured signal recovery problems where the prior imposed has a combinatorial flavor (see the next section and [13,14] for surveys).

Motivated by the recent advances of deep generative networks in learning complex data structures, in this paper we study the spiked random matrix models (1) and (2), where the planted signal  $y_*$  has a *generative network prior*. We assume that a generative neural network  $G : \mathbb{R}^k \rightarrow \mathbb{R}^n$  with  $k < n$ , has been trained on a data set of spikes, and the unknown spike  $y_* \in \mathbb{R}^n$  lies on the range of  $G$ , that is, we can write  $y_* = G(x_*)$  for some  $x_* \in \mathbb{R}^k$ . As a mathematical model for the trained  $G$ , we consider a network of the form:

$$G(x) = \text{relu}(W_d \dots \text{relu}(W_2 \text{relu}(W_1 x)) \dots), \quad (3)$$

with weight matrices  $W_i \in \mathbb{R}^{n_i \times n_{i-1}}$  and  $\text{relu}(x) = \max(x, 0)$  is applied entrywise. We furthermore assume that the network is expansive, that is,  $n = n_d > n_{d-1} > \dots > n_0 = k$ , and the weights have Gaussian entries. These modeling assumptions and their variants were used in [15–20].

Enforcing generative network priors has led to substantially fewer measurements needed for signal recovery than with traditional sparsity priors for a variety of signal recovery problems [17,21,22]. In the case of phase retrieval, [17,23] have shown that under the generative prior (3), efficient compressive phase retrieval is possible with sample complexity proportional (up to log factors) to the underlying signal dimensionality  $k$ . In contrast, for a sparsity-based prior, the best known polynomial time algorithms (convex methods [24–26], iterative thresholding [27–29], etc.) require a sample complexity proportional to the square of the sparsity level for stable recovery. Given that generative priors lead to no computational-statistical gap with compressive phase retrieval, one might anticipate that they will close other computational-statistical gaps as well. Indeed, [30] analyzed the spiked models (1) and (2) under a generative network prior similar to (3) and observed no computational-statistical gap in the asymptotic limit  $k, n, N \rightarrow \infty$  with  $n/k = \mathcal{O}(1)$  and  $N/n = \mathcal{O}(1)$ . For more details on this work and on the comparison of sparsity and generative priors, see Section 2.2.

### Our Contribution

In this paper we analyze the spiked matrix models (1) and (2) under a generative network prior in the nonasymptotic, finite data regime. We consider a  $d$ -layer feedforward generative network  $G : \mathbb{R}^k \rightarrow \mathbb{R}^n$  with architecture (3). We furthermore assume that the planted spike  $y_* \in \mathbb{R}^n$  lies on the range of  $G$ , that is, there exists a latent vector  $x_* \in \mathbb{R}^k$  such that  $y_* = G(x_*)$ .

To estimate  $y_*$ , we first find an estimate  $\hat{x}$  of the latent vector  $x_*$  and then use  $G(\hat{x})$  to estimate  $y_*$ . We thus consider the following minimization problem (under the conditions on the generative network specified below, it was shown in [15] that  $G$  is invertible and there exists a unique  $x_*$  that satisfies  $y_* = G(x)$ ):

$$\min_{x \in \mathbb{R}^k} f(x) := \frac{1}{4} \|G(x)G(x)^\top - M\|_F^2, \quad (4)$$

where:

- for the **Wishart model** (1) we take  $M = \Sigma_N - \sigma^2 I_n$  with  $\Sigma_N = Y^T Y / N$
- for the **Wigner model** (2) we take  $M = Y$ .

Despite the non-convexity and non-smoothness of the problem, our preliminary work in [31] shows that when the generative network  $G$  is expansive and has Gaussian weights, (4) enjoys a favorable optimization geometry. Specifically, every nonzero point outside two small neighborhoods around  $x_*$  and a negative multiple of it, has a descent direction which is given a.e. by the gradient of  $f$ . Furthermore, in [31] it is shown that the global minimum of  $f$  lies in the neighborhoods around  $x_*$  and has optimal reconstruction error. This result suggests that a first order optimization algorithm can succeed in efficiently solving (4), and no statistical-computational gap is present for the spiked matrix models with a (random) generative network prior in the finite data regime. In the current paper, we prove this conjecture by providing a polynomial-time subgradient method that minimizes the non-convex problem (4) and obtains information-theoretically optimal error rates.

Our main contribution can be summarized as follows. We analyze a subgradient method (Algorithm 1) for the minimization of (4) and show that after a polynomial number of steps  $\tilde{T}$  and up to polynomial factors in the depth  $d$  of the network, the iterate  $x_{\tilde{T}}$  satisfies the following reconstruction errors:

- in the **Spiked Wishart Model** :

$$\|G(x_{\tilde{T}}) - y_*\|_2 \lesssim \left(1 + \frac{\sigma^2}{\|y_*\|_2^2}\right) \sqrt{\frac{k \log(n)}{N}} \|y_*\|_2 \quad (5)$$

in the regime  $N \gtrsim k \log(n)$ ;

- in the **Spiked Wigner Model**:

$$\|G(x_{\tilde{T}}) - y_*\|_2 \lesssim \frac{v}{\|y_*\|_2^2} \sqrt{\frac{k \log(n)}{n}} \|y_*\|_2. \quad (6)$$

We notice that these bounds are information-theoretically optimal up to the log factors in  $n$ , and correspond to the best achievable in the case of a  $k$ -dimensional subspace prior. In particular, they imply that efficient recovery in the Wishart model is possible with a number of samples  $N$  proportional to the intrinsic dimension of the signal  $y_*$ . Similarly, the bound in the Spiked Wigner Model implies that imposing a generative network prior leads to a reduction of the noise by a factor of  $k/n$ .

---

#### Algorithm 1: Subgradient method for the minimization problem (4)

---

**Input:** Weights  $W_i$ , observation matrix  $M$ , step size  $\mu > 0$ , initial point

$x_0 \in \mathbb{R}^k \setminus \{0\}$  ;

```

1 for  $i = 0, 1, \dots$ , do
2   if  $f(x_i) > f(-x_i)$  then
3      $\tilde{x}_i \leftarrow -x_i$ ;
4   else
5      $\tilde{x}_i \leftarrow x_i$ 
6   Take  $v_{\tilde{x}_i} \in \partial f(\tilde{x}_i)$ ;
7    $x_{i+1} \leftarrow \tilde{x}_i - \mu v_{\tilde{x}_i}$  ;
```

**Output:**  $G(x_i), x_i$

---

## 2. Related Work

### 2.1. Sparse PCA and Other Computational-Statistical Gaps

A canonical problem in Statistics is finding the directions that explain most of the variance in a given cloud of data, and it is classically solved by Principal Component Analysis. Spiked covariance models were introduced in [1] to study the statistical performance

of this algorithm in the high dimensional regime. Under a spiked covariance model it is assumed that the data are of the form:

$$y_i = u_i y_* + \sigma z_i, \quad (7)$$

where  $\sigma > 0$ ,  $u_i \sim \mathcal{N}(0, 1)$  and  $z_i \sim \mathcal{N}(0, I_n)$  are independent and identically distributed, and  $y_*$  is the unit-norm planted spike. Each  $y_i$  is an i.i.d. sample from a centered Gaussian  $\mathcal{N}(0, \Sigma)$  with spiked covariance matrix given by  $\Sigma = y_* y_*^\top + \sigma^2 I_n$ , with  $y_*$  being the direction that explains most of the variance. The estimate of  $y_*$  provided by PCA is then given by the leading eigenvector  $\hat{y}$  of the empirical covariance matrix  $\Sigma_N = \frac{1}{N} \sum_{i=1}^N y_i y_i^\top$ , and standard techniques from high dimensional probability can be used to show that (we write  $f(n) \gtrsim g(n)$  if  $f(n) \geq Cn$  for some constant  $C > 0$  that might depend  $\sigma$  and  $\|y_*\|^2$ ). Similarly for  $f(n) \lesssim g(n)$  as long as  $N \gtrsim n$ ,

$$\min_{\epsilon = \pm} \|\epsilon \hat{y} - y_*\|_2 \lesssim \sqrt{\frac{n}{N}}, \quad (8)$$

with overwhelming probability. Note incidentally that the data matrix  $Y \in \mathbb{R}^{N \times n}$  with rows  $\{y_i^\top\}_i$  can be written as (1).

Bounds of the form (8), however, become uninformative in modern high dimensional regimes where the ambient dimension of the data  $n$  is much larger than, or on the order of, the number of samples  $N$ . Even worse, in the asymptotic regime  $n/N \rightarrow c > 0$  and for  $\sigma^2$  large enough, the spike  $y_*$  and the estimate  $\hat{y}$  become orthogonal [32], and minimax techniques show that no other estimators based solely on the data (7), can achieve better overlap with  $y_*$  [33].

In order to obtain consistent estimates and lower the sample complexity of the problem, therefore, additional prior information on the spike  $y_*$  has to be enforced. For this reason, in recent years various priors have been analyzed such as positivity [34], cone constraints [35] and sparsity [32,36]. In the latter case  $y_*$  is assumed to be  $s$ -sparse, and it can be shown (e.g., [33]) that for  $N \gtrsim s \log n$  and  $n \gtrsim s$ , the  $s$ -sparse largest eigenvector  $\hat{y}_s$  of  $\Sigma_N$

$$\hat{y}_s = \operatorname{argmax}_{y \in \mathcal{S}_2^{n-1}, \|y\|_0 \leq s} y^\top \Sigma_N y$$

satisfies with high probability the condition:

$$\min_{\epsilon = \pm} \|\epsilon \hat{y}_s - y_*\|_2 \lesssim \sqrt{\frac{s \log n}{N}}.$$

This implies, in particular, that the signal  $y_*$  can be estimated with a number of samples that scales linearly with its intrinsic dimension  $s$ . These rates are also minimax optimal; see for example [4] for the mean squared error and [2] for the support recovery. Despite these encouraging results, no currently known polynomial time algorithm achieves such optimal error rates and, for example, the covariance thresholding algorithm of [37] requires  $N \gtrsim s^2$  samples in order to obtain exact support recovery or estimation rate

$$\min_{\epsilon = \pm} \|\epsilon \hat{y}_s - y_*\|_2 \lesssim \sqrt{\frac{s^2 \log n}{N}},$$

as shown in [3]. In summary, only computationally intractable algorithms are known to reach the statistical limit  $N = \Omega(s)$  for Sparse PCA, while polynomial time methods are only sub-optimal, requiring  $N = \Omega(s^2)$ . Notably, [38] provided a reduction of Sparse PCA to the planted clique problem which is conjectured to be computationally hard.

Further strong evidence for the hardness of sparse PCA have been given in a series of recent works [39–43]. Other computational-statistical gaps have also been found and studied in a variety of other contexts such as sparse Gaussian mixture models [44], tensor principal component analysis [45], community detection [46] and synchronization over groups [47]. These works fit in the growing and important body of literature aiming at

understanding the trade-offs between statistical accuracy and computational efficiency in statistical inverse problems.

We finally note that many of the above mentioned problems can be phrased as recovery of a spike vector from a spiked random matrix. The difficulty can be viewed as arising from simultaneously imposing low-rankness and additional prior information on the signal (sparsity in case of Sparse PCA). This difficulty can be found in sparse phase retrieval as well. For example, [25] has shown that  $m = \mathcal{O}(s \log n)$  number of quadratic measurements are sufficient to ensure well-posedness of the estimation of an  $s$ -sparse signal of dimension  $n$  lifted to a rank-one matrix, while  $m \geq \mathcal{O}(s^2 / \log^2 n)$  measurements are necessary for the success of natural convex relaxations of the problem. Similarly, [48] studied the recovery of simultaneously low-rank and sparse matrices, showing the existence of a gap between what can be achieved with convex and tractable relaxations and nonconvex and intractable methods.

## 2.2. Inverse Problems with Generative Network Priors

Recently, in the wake of successes of deep learning, generative networks have gained popularity as a novel approach for encoding and enforcing priors in signal recovery problems. In one deep-learning-based approach, a dataset of “natural signals” is used to train a generative network in an unsupervised manner. The range of this network defines a low-dimensional set which, if successfully trained, contains or approximately contains, target signals of interest [19,21]. Non-convex optimization methods are then used for recovery by optimizing over the range of the network. We notice that allowing the algorithms the complete knowledge of the generative network architecture and of the learned weights is roughly analogous to allowing sparsity-based algorithms the knowledge of the basis or frame in which the signal is modeled as sparse.

The use of generative network for signal recovery has been successfully demonstrated in a variety of settings such as compressed sensing [21,49,50], denoising [16,51], blind deconvolution [22], inpainting [52] and many more [53–56]. In these papers, generative networks significantly outperform sparsity based priors at signal reconstruction in the low-measurement regime. This fundamentally leverages the fact that a natural signal can be represented more concisely by a generative network than by a sparsity prior under an appropriate basis. This characteristic has been observed even in untrained generative networks where the prior information is encoded only in the network architecture and has been used to devise state-of-the-art signal recovery methods [57–59].

Parallel to these empirical successes, a recent line of works have investigated theoretical guarantees for various statistical estimation tasks with generative network priors. Following the work of [15,21] gave global guarantees for compressed sensing, followed then by many others for various inverse problems [19,20,50,51,55]. In particular, in [17] the authors have shown that  $m = \Omega(k \log n)$  number of measurements are sufficient to recover a signal from random phaseless observations, assuming that the signal lies on the range of a generative network with latent dimension  $k$ . The same authors have then provided in [23] a polynomial time algorithm for recovery under the previous settings. Note that, contrary to the sparse phase retrieval problem, generative priors for phase retrieval allow for efficient algorithms with optimal sample complexity, up to logarithmic factors, with respect to the intrinsic dimension of the signal.

Further theoretical advances in signal recovery with generative network priors have been spurred by using techniques from statistical physics. Recently, [30] analyzed the spiked matrix models (1) and (2) with  $y_*$  in the range of a generative network with random weights, in the asymptotic limit  $k, n, N \rightarrow \infty$  with  $n/k = \mathcal{O}(1)$  and  $N/n = \mathcal{O}(1)$ . The analysis is carried out mainly for networks with sign or linear activation functions in the Bayesian setting where the latent vector is drawn from a separable distribution. The authors of [30] provide an Approximate Message Passing and a spectral algorithm, and they numerically observe no statistical-computational gap as these polynomial time methods are able to asymptotically match the information-theoretic optimum. In this asymptotic

regime, [60] further provided precise statistical and algorithmic thresholds for compressed sensing and phase retrieval.

### 3. Algorithm and Main Result

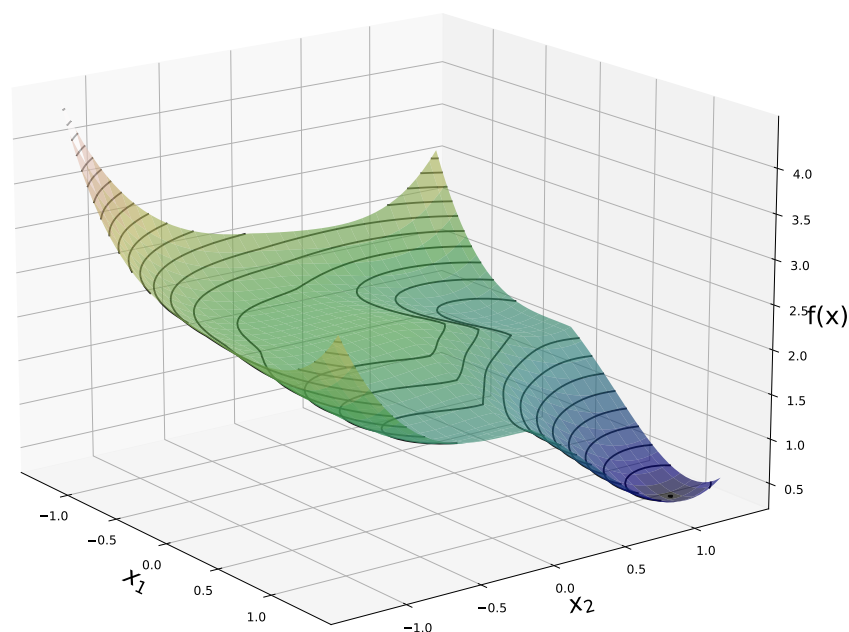
In this section we present an efficient and statistically-optimal algorithm for the estimation of the signal  $y_*$  given a spiked matrix  $Y$  of the form (1) or (2). The recovery method is detailed in Algorithm 1, and it is based on the direct optimization of the nonlinear least squares problem (4).

Applied in [16] for denoising and compressed sensing under generative network priors, and later used in [23] for phase retrieval, the first order optimization method described in Algorithm 1 leverages the theory of *Clarke subdifferentials* (the reader is referred to [61] for more details). As the objective function  $f$  is continuous and piecewise smooth, at every point  $x \in \mathbb{R}^k$  it has a Clarke subdifferential given by

$$\partial f(x) = \text{conv}\{v_1, v_2, \dots, v_T\}, \quad (9)$$

where  $\text{conv}$  denotes the convex hull of the vectors  $v_1, \dots, v_T$ , which are respectively the gradient of the  $T$  smooth functions adjoint at  $x$ . The vectors  $v_x \in \partial f(x)$  are the *subgradients* of  $f$  at  $x$ , and at a point  $x$  where  $f$  is differentiable it holds that  $\partial f(x) = \{\nabla f(x)\}$ .

The reconstruction method presented in Algorithm 1 is motivated by the landscape analysis of the minimization problem (4) for a network  $G$  with sufficiently-expansive Gaussian weights matrices. Under this assumption, we showed in [31] that (4) has a benign optimization geometry and in particular that for any nonzero point outside a neighborhood of  $x_*$  and a negative multiple of it, any subgradient of  $f$  is a direction of strict descent. Furthermore we showed that the points in the vicinity of the spurious negative multiple of  $x_*$  have function values strictly larger than those close to  $x_*$ . Figure 1 shows the expected value of  $f$  in the noiseless case,  $\nu = 0$  and  $N \rightarrow \infty$ , for a generative network with latent dimension  $k = 2$ . This plot highlights the global minimum at  $x_* = [1, 1]$ , and the flat region in near a negative multiple of  $x_*$ .



**Figure 1.** Expected value, with respect to the weights, of the objective function  $f$  in (4) in the noiseless case (see (16) for explicit formula), for a network with latent dimension  $k = 2$  and  $x_* = [1, 1]$ .

At each step, the subgradient method in Algorithm 1 checks if the current iterate  $x_i$  has a larger loss value than its negative multiple, and if so negates  $x_i$ . As we show in the proof



of our main result, this step will ensure that the algorithm will avoid the neighborhood around the spurious negative multiple of  $x_*$  and will converge to the neighborhood around  $x_*$  in a polynomial number of steps.

Below we make the following assumptions on the weight matrices of  $G$ .

**Assumption 1.** *The generative network  $G$  defined in (3), has weights  $W_i \in \mathbb{R}^{n_i \times n_{i-1}}$  with i.i.d. entries from  $\mathcal{N}(0, 1/n_i)$  and satisfying the expansivity condition with constant  $\epsilon > 0$ :*

$$n_{i+1} \geq c n_i \epsilon^{-2} \log(1/\epsilon) \tag{10}$$

for all  $i$  and a universal constant  $c > 0$ .

We note that in [31] the expansivity condition was more stringent, requiring an additional log factor. Since the publication of our paper, [62] has shown that the more relaxed assumption (10) suffices for ensuring a benign optimization geometry. Under Assumption 1, our main theorem below shows that the subgradient method in Algorithm 1 can estimate the spike  $y_*$  with optimal sample complexity and in a polynomial number of steps.

**Theorem 1.** *Let  $x_* \in \mathbb{R}^k$  nonzero and  $y_* = G(x_*)$  where  $G$  is a generative network satisfying Assumption 1 with  $\epsilon \leq K_1/d^{96}$ . Consider the minimization problem (4) and assume that the noise level  $\omega$  satisfies  $\omega \leq K_2 \|x_*\|_2^2 2^{-d}/d^{44}$  where:*

- for the **Spiked Wishart Model (1)** take  $M = \Sigma_N - \sigma^2 I_n$ , and

$$\omega := (\|y_*\|_2^2 + \sigma^2) \max \left\{ \sqrt{\frac{338k \log(3 n_1^d n_2^{d-1} \dots n_{d-1}^2 n)}{N}}, \frac{156k \log(3 n_1^d n_2^{d-1} \dots n_{d-1}^2 n)}{N} \right\};$$

- for the **Spiked Wigner Model (2)** take  $M = Y$ , and

$$\omega := \nu \sqrt{\frac{169k \log(3 n_1^d n_2^{d-1} \dots n_{d-1}^2 n)}{n}}.$$

Consider Algorithm 1 with nonzero  $x_0$  and  $\|x_0\|_2 < R_*$  where  $R_* \geq 5\|x_*\|_2/(2\sqrt{2})$ , and stepsize  $\mu = 2^{2d} K_3/(8d^4 R_*^2)$ . Then with probability at least  $1 - 2e^{-k \log n} - \sum_{i=1}^d e^{-Cn_{i-1}}$ ,  $0 < \|x_i\|_2 < R_*$  for any  $i \geq 1$ , there exists an integer  $T \leq K_4 f(x_0) 2^{2d}/(R_*^4 d^4 \epsilon)$  such that for any  $i \geq T$ :

$$\|x_{i+1} - x_*\|_2 \leq \rho_1^{i+1-T} \|x_T - x_*\|_2 + \rho_2 \frac{2^d}{\|x_*\|_2} \omega \tag{11}$$

$$\|G(x_{i+1}) - y_*\|_2 \leq \frac{1.2}{2^{d/2}} \rho_1^{i+1-T} \|x_T - x_*\|_2 + 1.3 \rho_2 \frac{\omega}{\|y_*\|_2} \tag{12}$$

where  $C > 0, K_1, \dots, K_4 > 0, \rho_1 \in (0, 1)$  and  $\rho_2 > 0$  are universal constants.

Note that the quantity  $2^{2d}$  in the hypotheses and conclusions of the theorem is an artifact of the scaling of the network and it should not be taken as requiring exponentially small noise or number of steps. Indeed under Assumption 1, the ReLU activation zeros out roughly half of the entries of its argument leading to an “effective” operator norm of approximately  $1/2$ . We furthermore notice that the dependence of the depth  $d$  is likely quite conservative and it was not optimized in the proof as the main objective was to obtain tight dependence on the intrinsic dimension of the signal  $k$ . As shown in the numerical experiments, the actual dependence on the depth is much better in practice. Finally, observe that despite the nonconvex nature of the objective function in (4) we obtain

a rate of convergence which is not directly dependent on the dimension of the signal, reminiscent of what happens in the convex case.

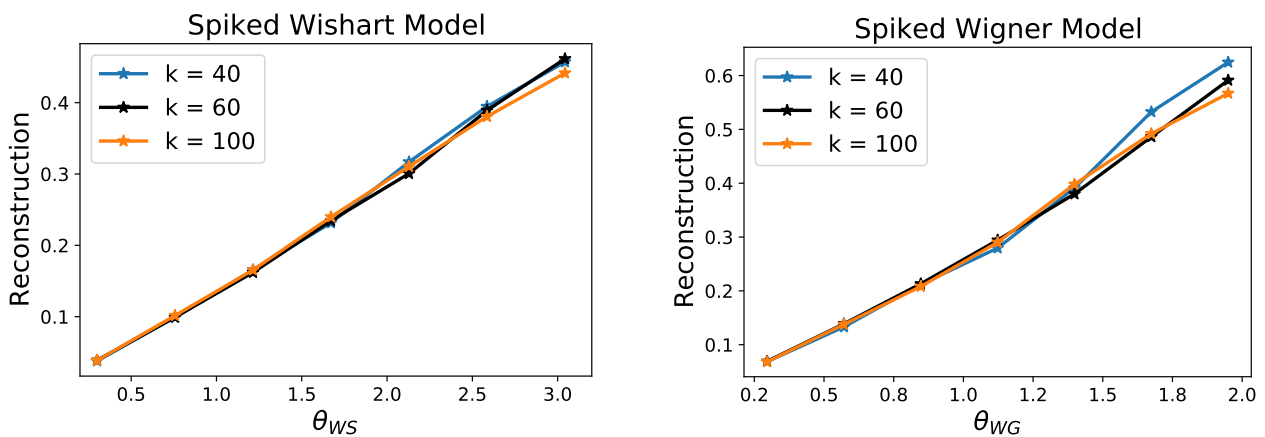
The quantity  $\omega$  in Theorem 1 can be interpreted as the intrinsic noise level of the problem (inverse SNR). The theorem guarantees that in a polynomial number of steps the iterates of the subgradient method will converge to  $x_*$  up to  $\omega$ . For  $\tilde{T}$  large enough  $G(x_{\tilde{T}})$  will satisfy the rate-optimal error bounds (5) and (6).

*Numerical Experiments*

We illustrate the predictions of our theory by providing results of Algorithm 1 on a set of synthetic experiments. We consider 2-layer generative networks with ReLU activation functions, hidden layer of dimension  $n_1 = 500$ , output dimension  $n_2 = n = 1500$  and varying number of latent dimension  $k \in [40, 60, 100]$ . We randomly sample the weights of the matrix independently from  $\mathcal{N}(0, 2/n_i)$  (this scaling removes that  $2^d$  dependence in Theorem 1). We then consider data  $Y$  according the spiked models (1) and (2), where  $x_* \in \mathbb{R}^k$  is chosen so that  $y_* = G(x_*)$  has unit norm. For the Wishart model, we vary the number of samples  $N$ ; and for the Wigner model, we vary the noise level  $\nu$  so that the following quantities remain constant for the different networks with latent dimension  $k$ :

$$\theta_{ws} := \sqrt{k \log(n_1^2 n) / N}, \quad \theta_{wg} := \nu \sqrt{k \log(n_1^2 n) / n}.$$

In Figure 2 we plot the reconstruction error given by  $\|G(x) - y_*\|_2$  against  $\theta_{ws}$  and  $\theta_{wg}$ . As predicted by Theorem 1, the errors scale linearly with respect to these control parameters, and moreover the overlap of these plots confirms that these rates are tight with respect to the order of  $k$ .



**Figure 2.** Reconstruction error for the recovery of a spike  $y_* = G(x_*)$  in the Wishart and Wigner models with random generative network priors. Each point corresponds to the average over 50 random drawing of the network weights and samples. These plots demonstrate that the reconstruction errors follow the scalings established by Theorem 1.

**4. Recovery Under Deterministic Conditions**

We will derive Theorem 1 from Theorem 3, below, which is based on a set of deterministic conditions on the weights of the matrix and the noise. Specifically, we consider the minimization problem (4) with

$$M = G(x_*)G(x_*)^T + H$$

for an unknown symmetric matrix  $H \in \mathbb{R}^{n \times n}$ , nonzero  $x_* \in \mathbb{R}^k$ , and a given  $d$ -layer feed forward generative network  $G$  as in (3).

In order to describe the main deterministic conditions on the generative network  $G$ , we begin by introducing some notation. For  $W \in \mathbb{R}^{n \times k}$  and  $x \in \mathbb{R}^k$ , we define the operator



$W_{+,x} := \text{diag}(Wx > 0)W$  such that  $\text{relu}(Wx) = W_{+,x}x$ . Moreover, we let  $W_{1,+,x} = (W_1)_{+,x} = \text{diag}(W_1x > 0)W_1$ , and for  $2 \leq i \leq d$  we define recursively

$$W_{i,+,x} = \text{diag}(W_i, \Pi_{j=i-1}^1 W_{j,+,x} > 0)W_i,$$

where  $\Pi_{i=d}^1 W_i = W_d W_{d-1} \dots W_1$ . Finally we let  $\Lambda_x = \Pi_{j=d}^1 W_{j,+,x}$  and note that  $G(x) = \Lambda_x x$ . With this notation we next recall the following deterministic condition on the layers of the generative network.

**Definition 2** (Weight Distribution Condition [15]). *We say that  $W \in \mathbb{R}^{n \times k}$  satisfies the Weight Distribution Condition (WDC) with constant  $\epsilon > 0$  if for all nonzero  $x_1, x_2 \in \mathbb{R}^k$ :*

$$\|W_{+,x_1}^\top W_{+,x_2} - Q_{x_1,x_2}\|_2 \leq \epsilon,$$

where

$$Q_{x_1,x_2} = \frac{\pi - \theta_{x_1,x_2}}{2\pi} I_k + \frac{\sin \theta_{x_1,x_2}}{2\pi} M_{\hat{x}_2 \leftrightarrow \hat{x}_1}$$

and  $\theta_{x_1,x_2} = \angle(x_1, x_2)$ ,  $\hat{x}_1 = x_1 / \|x_1\|_2$ ,  $\hat{x}_2 = x_2 / \|x_2\|_2$ ,  $I_k$  is the  $k \times k$  identity matrix and  $M_{\hat{x}_1 \leftrightarrow \hat{x}_2}$  is the matrix that sends  $\hat{x}_1 \mapsto \hat{x}_2$ ,  $\hat{x}_2 \mapsto \hat{x}_1$ , and with kernel  $\text{span}(\{x_1, x_2\})^\perp$ .

Note that  $Q_{x_1,x_2}$  is the expected value of  $W_{+,x_1}^\top W_{+,x_2}$  when  $W$  has rows  $w_i \sim \mathcal{N}(0, I_k/n)$ , and if  $x_1 = x_2$  then  $Q_{x_1,x_2}$  is an isometry up to the scaling factor  $1/2$ . Below we will say that a  $d$ -layer generative network  $G$  of the form (3), satisfies the WDC with constant  $\epsilon > 0$  if every weight matrix  $W_i$  has the WDC with constant  $\epsilon$  for all  $i = 1, \dots, d$ .

The WDC was originally introduced in [15], and ensures that the angle between two vectors in the latent space is approximately preserved at the output layer and, in turn, it guarantees the invertibility of the network. Assumption 1 will guarantee that the generative network  $G$  satisfies the WDC with high probability.

We are now able to state our recovery guarantees for a spike  $y_*$  under deterministic conditions on the network  $G$  and noise  $H$ .

**Theorem 3.** *Let  $d \geq 2$  and assume the generative network (3) has weights  $W_i \in \mathbb{R}^{n_i \times n_{i-1}}$  satisfying the WDC with constant  $0 < \epsilon \leq K_1/d^{96}$ . Consider Algorithm 1 with  $M = G(x_*)G(x_*)^\top + H$ ,  $x_* \in \mathbb{R}^k \setminus \{0\}$  and  $H$  a symmetric matrix satisfying:*

$$\|\Lambda_x^\top H \Lambda_x\|_2 \leq \frac{\omega}{2^d}, \quad \text{and} \quad \omega \leq K_2 \frac{\|x_*\|_2^{2-d}}{d^{44}}. \tag{13}$$

Take  $x_0$  nonzero and with  $\|x_0\|_2 < R_*$  where  $R_* \geq 5\|x_*\|_2 / (2\sqrt{2})$ ,  $\mu = 2^{2d}K_3 / (8d^4R_*^2)$ . Then the iterates  $\{x_i\}_{i \geq 0}$  generated by the Algorithm 1 satisfy  $0 < \|x_i\|_2 < R_*$  and obey to the following:

(A) *there exists an integer  $T \leq K_4 \frac{f(x_0)2^{2d}}{R_*^4 d^4 \epsilon}$  such that*

$$\|x_T - x_*\|_2 \leq K_5 d^{14} \sqrt{\epsilon} \|x_*\|_2 + K_6 2^d d^{10} \omega \|x_*\|_2^{-1}$$

(B) *for any  $i \geq T$ :*

$$\|x_{i+1} - x_*\|_2 \leq \rho_1^{i+1-T} \|x_T - x_*\|_2 + \rho_2 \frac{2^d}{\|x_*\|_2} \omega \tag{14}$$

$$\|G(x_{i+1}) - G(x_*)\|_2 \leq \frac{1.2}{2^{d/2}} \rho_1^{i+1-T} \|x_T - x_*\|_2 + 1.3 \rho_2 \frac{\omega}{\|y_*\|_2} \tag{15}$$

where  $K_1, \dots, K_6 > 0$ ,  $\rho_1 \in (0, 1)$  and  $\rho_2 > 0$  are universal constants.

Theorem 1 follows then from Theorem 3 after proving that with high probability the

spectral norm of  $\Lambda_x^\top H \Lambda_x$ , where  $H = M - y_* y_*^\top$ , can be upper bounded by  $\omega/2^d$ , and the weights of the network  $G$  satisfy with high probability the WDC.

In the rest of this section we will describe the main steps and tools needed to prove Theorem 3.

#### 4.1. Technical Tools and Outline of the Proofs

Our proof strategy for Theorem 3 can be summarized as follows:

1. In Proposition A1 (Appendix A.3) we show that the iterates  $\{x_i\}_{i=1}$  of the Algorithm 1 stay inside the Euclidean ball of radius  $R_*$  and remain nonzero for all  $i \geq 1$ .
2. We then identify two small Euclidean balls  $\mathcal{B}_+$  and  $\mathcal{B}_-$  around respectively  $x_*$  and  $-\rho_d x_*$ , where  $\rho_d \in (0, 1)$  only depends on the depth of the network. In Proposition A2 we show that after a polynomial number of steps, the iterates  $\{x_i\}$  of the Algorithm 1 enter the region  $\mathcal{B}_+ \cup \mathcal{B}_-$  (Appendix A.4).
3. We show, in Proposition A3, that the negation step causes the iterates of the algorithm to avoid the spurious point  $-\rho_d x_*$  and actually enter  $\mathcal{B}_+$  within a polynomial number of steps (Appendix A.5).
4. We finally show in Proposition A4, that in  $\mathcal{B}_+$  the loss function  $f$  enjoys a favorable convexity-like property, which implies that the iterates  $\{x_i\}$  will remain in  $\mathcal{B}_+$  and eventually converge to  $x_*$  up to the noise level (Appendix A.6).

One of the main difficulties in the analysis of a subgradient method in Algorithm 1 is the lack of smoothness of the loss function  $f$ . We show that the WDC allows us to overcome this issue by showing that the subgradients of  $f$  are uniformly close, up to the noise level, to the vector field  $h_x \in \mathbb{R}^k$ :

$$h_x := \frac{1}{2^{2d}} (xx^\top - \tilde{h}_x \tilde{h}_x^\top) x,$$

where  $\tilde{h}_x$  is continuous for nonzero  $x$  (see Appendix A.2). We show furthermore that  $h_x$  is locally Lipschitz, which allows us to conclude that the gradient method decreases the value of the loss function until eventually reaching  $\mathcal{B}_+ \cup \mathcal{B}_-$  (Appendix A.4).

Using the WDC, we show that the loss function  $f$  is uniformly close to

$$f_E(x) = \frac{1}{2^{2d+2}} \left( \|x\|_2^4 + \|x_*\|_2^4 - 2\langle x, \tilde{h}_x \rangle^2 \right). \quad (16)$$

A direct analysis of  $f_E$  reveals that its values inside  $\mathcal{B}_-$  are strictly larger than those inside  $\mathcal{B}_+$ . This property extends to  $f$  as well, and guarantees that the gradient method will not converge to the spurious point  $-\rho_d x_*$  (Appendix A.5).

**Author Contributions:** Conceptualization, P.H. and V.V.; Formal analysis, J.C., Writing—original draft, J.C.; Writing - review & editing J.C. and P.H.; Supervision P.H. and V.V. All authors have read and agreed to the published version of the manuscript.

**Funding:** PH was partially supported by NSF CAREER Grant DMS-1848087 and NSF Grant DMS-2022205.

**Conflicts of Interest:** The authors declare no conflict of interest.

### Appendix A. Supporting Lemmas and Proof of Theorem 3

The proof of Theorem 3 is provided in Appendix A.7. We begin this section with a set of preliminary results and supporting lemmas.

#### Appendix A.1. Notation

We collect the notation that is used throughout the paper. For any real number  $a$ , let  $\text{relu}(a) = \max(a, 0)$  and for any vector  $v \in \mathbb{R}^n$ , denote the entrywise application of  $\text{relu}$  as  $\text{relu}(v)$ . Let  $\text{diag}(Wx > 0)$  be the diagonal matrix with  $i$ -th diagonal element equal to 1 if  $(Wx)_i > 0$  and 0 otherwise. For any vector  $x$  we denote with  $\|x\|$  its Euclidean

norm and for any matrix  $A$  we denote with  $\|A\|$  its spectral norm and with  $\|A\|_F$  its Frobenius norm. The euclidean inner product between two vectors  $a$  and  $b$  is  $\langle a, b \rangle$ , while for two matrices  $A$  and  $B$  their Frobenius inner product will be denoted by  $\langle A, B \rangle_F$ . For any nonzero vector  $x \in \mathbb{R}^n$ , let  $\hat{x} = x/\|x\|$ . For a set  $S$  we will write  $|S|$  for its cardinality and  $S^c$  for its complement. Let  $\mathcal{B}(x, r)$  be the Euclidean ball of radius  $r$  centered at  $x$ , and  $S^{k-1}$  be the unit sphere in  $\mathbb{R}^k$ . Let  $\theta_0 = \angle(x, x_*)$  and for  $i \geq 0$  let  $\theta_{i+1} = g(\theta_i)$  where  $g$  is defined in (A1). We will write  $\gamma = \Omega(\delta)$  to mean that there exists a positive constant  $C$  such that  $\gamma \geq C\delta$  and similarly  $\gamma = \mathcal{O}(\delta)$  if  $\gamma \leq C\delta$ . Additionally we will use  $a = b + \mathcal{O}_1(\delta)$  when  $\|a - b\| \leq \delta$ , where the norm is understood to be the absolute value for scalars, the Euclidean norm for vectors and the spectral norm for matrices.

Appendix A.2. Preliminaries

For later convenience we will define the following vectors:

$$\begin{aligned} p_x &:= \Lambda_x^\top \Lambda_x x, \\ q_x &:= \Lambda_x^\top \Lambda_{x_*} x_*, \\ \bar{v}_x &:= [p_x p_x^\top - q_x q_x^\top] x, \\ \eta_x &:= \Lambda_x^\top H \Lambda_x x. \end{aligned}$$

Note then that when  $f$  is differentiable at  $x$ , then  $\bar{v}_x := \nabla f(x) = \bar{v}_x - \eta_x$  and in particular when  $H = 0$  then we have  $\bar{v}_x = \bar{v}_x$ .

The following function controls how the angles are contracted by a ReLU layer:

$$g(\theta) := \cos^{-1} \left( \frac{(\pi - \theta) \cos \theta + \sin \theta}{\pi} \right). \tag{A1}$$

As we mentioned in Section 4.1, our analysis is based on showing that the subgradients of  $f$  are uniformly close to the vector field  $h_x$  given by

$$h_x := \frac{1}{2^{2d}} (x x^\top - \tilde{h}_x \tilde{h}_x^\top) x, \tag{A2}$$

where

$$\tilde{h}_x := \left( \prod_{i=0}^{d-1} \frac{\pi - \bar{\theta}_i}{\pi} \right) x_* + \sum_{i=1}^{d-1} \frac{\sin \theta_i}{\pi} \left( \prod_{j=i+1}^{d-1} \frac{\pi - \bar{\theta}_j}{\pi} \right) \frac{\|x_*\|}{\|x\|} x, \tag{A3}$$

and  $\theta_i := g(\bar{\theta}_{i-1})$  for  $g$  given by (A1) and  $\theta_0 = \angle(x, y)$ .

**Lemma A1** (Lemma 8 in [15]). *Suppose that  $d \geq 2$  and the WDC holds with  $\epsilon < 1/(16\pi d^2)^2$ , then for all nonzero  $x, x_* \in \mathbb{R}^k$ ,*

$$\langle \Lambda_x x, \Lambda_{x_*} x_* \rangle \geq \frac{1}{4\pi} \frac{1}{2^d} \|x\|_2 \|x_*\|, \tag{A4}$$

$$\|\Lambda_x^\top \Lambda_{x_*} x_* - \frac{\tilde{h}_x}{2^d}\| \leq 24 \frac{d^3 \sqrt{\epsilon}}{2^d} \|x_*\|, \text{ and} \tag{A5}$$

$$\|\Lambda_x\|^2 \leq \frac{1}{2^d} (1 + 2\epsilon)^d \leq \frac{1 + 4\epsilon d}{2^d} \leq \frac{13}{12} \frac{1}{2^d}. \tag{A6}$$

**Proof.** The first two bounds can be found in [15] (Lemma 8). The third bound follows by noticing that the WDC implies:

$$\|\Lambda_x\|^2 \leq \prod_{i=d}^1 \|W_{i+,x}\|^2 \leq \frac{1}{2^d} (1 + 2\epsilon)^d \leq \frac{1 + 4\epsilon d}{2^d} \leq \frac{13}{12} \frac{1}{2^d},$$

where we used  $\log(1 + z) \leq z$  and  $e^z \leq (1 + 2z)$  for all  $0 \leq z \leq 1$ .  $\square$

The next lemma shows that the noiseless gradient  $\bar{v}_x$  concentrates around  $h_x$ .

**Lemma A2.** *Suppose  $d \geq 2$  and the WDC holds with  $\epsilon < 1/(16\pi d^2)^2$ , then for all nonzero  $x, x_\star \in \mathbb{R}^k$ :*

$$\|\bar{v}_x - h_x\| \leq 86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} \max(\|x_\star\|^2, \|x\|^2) \|x\|.$$

We now use the characterization of the Clarke subdifferential given in (9) to derive a bound on the concentration of  $v_x \in \partial f(x)$  around  $h_x$  up to the noise level.

**Lemma A3.** *Under the assumptions of Lemma A2, and assuming  $\|\Lambda_x^\top H \Lambda_x\| \leq \omega/2^d$ , for any  $v_x \in \partial f(x)$ :*

$$\|v_x - h_x\| \leq 86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} \max(\|x_\star\|^2, \|x\|^2) \|x\| + \frac{\omega}{2^d} \|x\|$$

From the above and the bound on the noise level  $\omega$  we can bound the norm of the step  $v_x$  in the Algorithm 1.

**Lemma A4.** *Under the assumptions of Lemma A3, and assuming that  $\omega$  satisfies (13), for any  $v_x \in \partial f(x)$ :*

$$\|v_x\| \leq \frac{4}{2^{2d}} d^2 \max(\|x\|^2, \|x_\star\|^2) \|x\|.$$

#### Appendix A.3. Iterates Stay Bounded

In this section we prove that all the iterates  $\{x_i\}$  generated by Algorithm 1 remain inside the Euclidean ball  $\mathcal{B}(0, R_\star)$  where  $R_\star \geq \sqrt{2} C_\star \|x_\star\|$  and  $C_\star = 5/4$ .

**Lemma A5.** *Let the assumptions of Lemma A4 be satisfied,  $\mu = K_3 2^{2d} / (8d^4 R_\star^2)$  with  $512\pi^2 K_3 < 1$ . Then for any  $x \in \mathbb{R}^k$  with  $C_\star \|x_\star\| < \|x\| \leq R_\star$  and any  $\lambda \in [0, 1]$ , it holds that  $\|x - \lambda \mu v_x\| \leq \|x\|$ .*

From the previous lemma we can now derive the boundedness of the iterates of the Algorithm 1.

**Proposition A1.** *Under the assumptions of Theorem 3, if  $x \in \mathcal{B}(0, R_\star) \setminus \{0\}$  it follows that  $x - \lambda \mu v_x \in \mathcal{B}(0, R_\star) \setminus \{0\}$ . Furthermore if  $x_0 \in \mathcal{B}(0, R_\star) \setminus \{0\}$ , the iterates  $\{x_i\}_{i=1}$  of the Algorithm 1 satisfy  $x_i - \lambda \mu v_{x_i} \in \mathcal{B}(0, R_\star) \setminus \{0\}$  for all  $i \geq 1$  and  $\lambda \in [0, 1]$ .*

**Proof.** Assume  $C_\star \|x_\star\| < \|x\| \leq R_\star$ , then the conclusions follows from Lemma A5. Assuming instead that  $\|x\| \leq C_\star \|x_\star\|$ , note that

$$\|x - \lambda \mu v_x\| \leq \|x\| + \mu \|v_x\| \leq (1 + \frac{1}{2^{2d}}) \|x\| \leq R_\star$$

using Lemma A4,  $d \geq 2$  and the assumptions on  $\mu$  and  $R_\star$ . Finally observe that if  $x_i \in \mathcal{B}(0, R_\star)$  then the same holds for  $\tilde{x}_i$ .

Finally if for some  $i \geq 0$  it was the case that  $0 = \tilde{x}_i - \lambda \mu v_{\tilde{x}_i}$ , then this would imply that  $\|\tilde{x}_i\| = \lambda \mu \|v_{\tilde{x}_i}\|$  which cannot happen because by Lemma A4 and the choice of the step size it holds that  $\mu \|v_{\tilde{x}_i}\| \leq \|\tilde{x}_i\|/8$ .  $\square$

#### Appendix A.4. Convergence to $\mathcal{S}_\beta$

We define the set  $\mathcal{S}_\beta$  outside which we can lower bound the norm of  $\|h_x\|$  as

$$\mathcal{S}_\beta := \{x \in \mathbb{R}^k \mid \|h_x\| \leq \frac{\beta}{2^{2d}} \max(\|x\|^2, \|x_\star\|^2) \|x\|\},$$

where we take

$$\beta = 7 \cdot (86d^4\sqrt{\epsilon} + 2^d\omega\|x_\star\|^{-2}). \tag{A7}$$

Outside the set  $\mathcal{S}_\beta$  the sub-gradients of  $f$  are bounded below and the landscape has favorable optimization geometry.

**Lemma A6.** *Let  $x \in \mathcal{B}(0, R_\star) \setminus \{0\} \cap \mathcal{S}_\beta^c$ , then for all  $v_x \in \partial f(x)$*

$$\|v_x\| \geq 6 \left( \frac{\max(\|x\|^2, \|x_\star\|^2)\|x\|}{2^{2d}} 86d^4\sqrt{\epsilon} + \frac{\omega}{2^d}\|x\| \right). \tag{A8}$$

Moreover let  $\lambda \in [0, 1]$  and  $x_\lambda = x - \lambda\mu v_{x_\lambda}$  then

$$\|v_x - v_{x_\lambda}\| \leq \frac{15}{16}\|v_x\|, \tag{A9}$$

for all  $v_x \in \partial f(x)$  and  $v_{x_\lambda} \in \partial f(x_\lambda)$ .

Based on the previous lemma we can prove the main result of this section.

**Proposition A2.** *Under the assumptions of Theorem 3, if  $x_i \in \mathcal{B}(0, R_\star) \setminus \{0\} \cap \mathcal{S}_\beta^c$  then*

$$f(x_{i+1}) - f(x_i) \leq -\frac{KR_\star^4 d^4 \epsilon}{2^{2d}} \tag{A10}$$

for some numerical constant  $K > 0$ . Moreover there exists an integer  $T \leq \frac{f(x_0)2^{2d}}{KR_\star^4 d^4 \epsilon}$  such  $x_{i+T} \in \mathcal{S}_\beta$ .

**Proof.** Let  $x_i \notin \mathcal{S}_\beta$  and assume that  $f(x_i) \leq f(-x_i)$ , then  $\tilde{x}_i = x_i$ . By the mean value theorem for Clarke subdifferentials [61] (Theorem 8.13), there exists  $\lambda \in (0, 1)$  such that for  $x_{i,\lambda} = x_i - \mu\lambda v_{x_i}$  and a  $v_{x_{i,\lambda}} \in \partial f(x_{i,\lambda})$  it holds that

$$\begin{aligned} f(x_{i+1}) - f(\tilde{x}_i) &= \langle v_{x_{i,\lambda}}, -\mu v_{x_i} \rangle \\ &= \langle v_{x_i}, -\mu v_{x_i} \rangle + \langle v_{x_{i,\lambda}} - v_{x_i}, -\mu v_{x_i} \rangle \\ &\leq -\mu\|v_{x_i}\|(\|v_{x_i}\| - \|v_{x_{i,\lambda}} - v_{x_i}\|) \\ &\leq -\frac{\mu}{16}\|v_{x_i}\|^2 \end{aligned} \tag{A11}$$

where the first inequality follows from the triangle inequality and the second from Equation (A9). Next observe that by (A8)

$$\|v_x\|^2 \geq 36 \frac{\max(\|x\|^2, \|x_\star\|^2)^2\|x\|^2}{2^{4d}} 86^2 d^8 \epsilon$$

which together with the definition of  $\mu$  and (A11) gives (A10).

Next take  $x_i \notin \mathcal{S}_\beta$  and assume  $f(x_i) > f(-x_i)$ , so that  $\tilde{x}_i = -x_i$ . Observe that

$$f(x_{i+1}) - f(x_i) < f(x_{i+1}) - f(-x_i) = f(x_{i+1}) - f(\tilde{x}_i),$$

we obtain then (A10) proceeding as before.

Finally the claim on the maximum number of iterations directly follows by a telescopic sum on (A10) and  $f(\cdot) \geq 0$ .  $\square$

#### Appendix A.5. Convergence to a Neighborhood Around $x_\star$

In the previous section we have shown that after a finite number of steps the iterates  $\{x_i\}$  of the Algorithm 1 will enter in the region  $\mathcal{S}_\beta$ . In this section we show that, thanks to the negation step in the descent algorithm, they will eventually be confined in a neighborhood of  $x_\star$ .

The following lemma shows that  $\mathcal{S}_\beta$  is contained inside two balls around  $x_\star$  and  $-x_\star$ .

**Lemma A7.** Suppose  $8\pi d^6 \sqrt{\beta} \leq 1$ , then we have  $\mathcal{S}_\beta \subset \mathcal{B}^+ \cup \mathcal{B}^-$  where

$$\mathcal{B}^+ := \mathcal{B}(x_*, R_1 \beta d^{10} \|x_*\|) \quad \text{and} \quad \mathcal{B}^- := \mathcal{B}(-\rho_d x_*, R_2 \sqrt{\beta} d^{10} \|x_*\|),$$

where  $R_1, R_2$  are numerical constants and  $0 < \rho_d < 1$  such that  $\rho_d \rightarrow 1$  as  $d \rightarrow \infty$ .

We furthermore observe that the ball around  $-x_*$  has values strictly higher than the one around  $x_*$ .

**Lemma A8.** Suppose that  $d \geq 2$ , the WDC holds with  $\epsilon < 1/(16\pi d^2)^2$  and  $H$  satisfies (13). Then for any  $\phi_d \in [\rho_d, 1]$ , it holds that

$$f(x) < f(y), \tag{A12}$$

for all  $x \in \mathcal{B}(\phi_d x_*, \epsilon \|x_*\| d^{-12})$  and  $y \in \mathcal{B}(-\phi_d x_*, \epsilon \|x_*\| d^{-12})$  where  $\epsilon < 1$  is a universal constant.

The main result of this section is about the convergence to a neighborhood of  $x_*$  of the iterates  $\{x_i\}$ .

**Proposition A3.** Under the assumptions of Theorem 3, if  $x_i \in \mathcal{B}(0, R_*) \setminus \{0\}$ , then there exists a finite number of steps  $T \leq K \frac{f(x_i) 2^{2d}}{R_*^4 d^4 \epsilon}$  such that  $x_{i+T} \in \mathcal{B}(x_*, R_1 \beta d^{10} \|x_*\|)$ . In particular it holds that

$$\|x_{i+T} - x_*\| \leq C_1 d^{14} \sqrt{\epsilon} \|x_*\| + C_2 2^d d^{10} \omega \|x_*\|^{-1}. \tag{A13}$$

**Proof.** Either  $x_i \in \mathcal{S}_\beta$  or by Proposition A2 there exist  $T'$  such that  $x_{i+T'} \in \mathcal{S}_\beta$ . By the choice of  $\epsilon > 0$ , the definition of  $\beta$  in (A7), and the assumption on the noise level (13), it follows that the hypotheses of Lemma A7 are satisfied. We therefore define the two neighborhoods  $\mathcal{S}_\beta^+ := \mathcal{S}_\beta \cap \mathcal{B}^+$  and  $\mathcal{S}_\beta^- := \mathcal{S}_\beta \cap \mathcal{B}^-$  and conclude that either  $x_{i+T'} \in \mathcal{S}_\beta^+$  or  $x_{i+T'} \in \mathcal{S}_\beta^-$ .

We next notice that  $\mathcal{S}_\beta^+ \subseteq \mathcal{B}(x_*, \epsilon \|x_*\| d^{-12})$  and  $\mathcal{S}_\beta^- \subseteq \mathcal{B}(-\rho_d x_*, \epsilon \|x_*\| d^{-12})$ . We can then use Lemma A8 to conclude that by the negation step, if  $x_{i+T'} \in \mathcal{S}_\beta^+$  then  $\tilde{x}_{i+T'} \in \mathcal{S}_\beta^+$ , otherwise we will have  $\tilde{x}_{i+T'} \in -\mathcal{S}_\beta^-$ .

We now analyze the case  $x_{i+T'} \in -\mathcal{S}_\beta^- \cap \mathcal{S}_\beta^c$ . Applying again Proposition A2, we have that there exists an integer  $T'' \leq K \frac{f(x_{i+T'}) 2^{2d}}{R_*^4 d^4 \epsilon}$  such that  $x_{i+T'+T''} \in \mathcal{S}_\beta$ . Furthermore Proposition A2 implies that  $f(x_{i+T'+T''}) < f(x_{i+T'})$ , while from Lemma A8 we know that  $f(x_{i+T'}) < f(y)$  for all  $y \in \mathcal{S}_\beta^-$ . We conclude therefore that  $x_{i+T'+T''}$  must be in  $\mathcal{S}_\beta^+$ .

In summary we obtained that there exists an integer  $T \leq K \frac{f(x_0) 2^{2d}}{R_*^4 d^4 \epsilon}$  such that  $x_{i+T} \in \mathcal{S}_\beta^+ \subset \mathcal{B}^+$ . Finally Equation (A13) follows from the definition of  $\beta$  in (A7) and  $\mathcal{B}^+$ .  $\square$

Appendix A.6. Convergence to  $x_*$  up to Noise

**Lemma A9.** Suppose the WDC holds with  $\epsilon < 1/(200^4 d^6)$ , then for all  $x \in \mathcal{B}(x_*, d\sqrt{\epsilon} \|x_*\|)$  and  $v_x \in \partial_x f(x)$ , it holds that

$$\|v_x - \frac{\tau_1}{2^{2d}} \|x_*\|^2 (x - x_*)\| \leq \tau_2 \frac{\|x_*\|^2}{2^{2d}} \|x - x_*\| + \tau_3 \frac{\omega}{2^d} \|x_*\|$$

where  $\tau_1 = 21/5$ ,  $\tau_2 = 17/5$  and  $\tau_3 = (1 + 1/(400)^2)$ .

Based on the previous condition on the direction of the sub-gradients we can then prove that the iterates of the algorithm converge to  $x_*$  up to noise.



**Proposition A4.** Under the assumptions of Theorem (3), if  $x_T \in \mathcal{B}^+$  then for any  $i \geq T$  it holds that  $x_i \in \mathcal{B}^+$  and furthermore

$$\|x_{i+1} - x_\star\| \leq \rho_1^{i+1-T} \|x_T - x_\star\| + \rho_2 \frac{2^d}{\|x_\star\|} \omega \tag{A14}$$

where  $\rho_1 \in (0, 1)$  and  $\rho_2 > 0$  are numerical constants.

**Proof.** If  $i = T$ , by Proposition A2, it follows that  $x_i = \tilde{x}_i \in \mathcal{B}^+$ . Furthermore, the assumptions of Lemma A9 are satisfied and we can write:

$$\begin{aligned} \|x_{i+1} - x_\star\| &= \|\tilde{x}_i - \mu v_{\tilde{x}_i} - x_\star\| \\ &\leq \left(1 - \frac{\mu \tau_1 \|x_\star\|^2}{2^{2d}}\right) \|\tilde{x}_i - x_\star\| + \mu \|v_{\tilde{x}_i} - \frac{\tau_1}{2^{2d}} x_\star\|^2 \|\tilde{x}_i - x_\star\| \\ &\leq \left(1 - \frac{\mu(\tau_1 - \tau_2) \|x_\star\|^2}{2^{2d}}\right) \|\tilde{x}_i - x_\star\| + \mu \tau_3 \frac{\omega}{2^d} \|x_\star\|. \end{aligned} \tag{A15}$$

Next recall that  $\omega$  satisfies (13),  $\mu = K_3 2^{2d} / (8d^4 R_\star^2)$  and  $R_\star^2 \geq 25 \|x_\star\|^2 / 8$ , then

$$\begin{aligned} \|x_{i+1} - x_\star\| &\leq \left(1 - \frac{\mu(\tau_1 - \tau_2) \|x_\star\|^2}{2^{2d}}\right) \|\tilde{x}_i - x_\star\| + \mu \frac{K_2 \tau_3}{2^{2d} d^{44}} \|x_\star\|^3 \\ &\leq \left[1 - (\tau_1 - \tau_2) \frac{K_3}{25d^4} + K_2 \frac{K_3}{25d^{48}}\right] \|x_\star\| \end{aligned}$$

Therefore for  $K_2$  and  $K_3$  small enough we obtain that  $x_{i+1} \in \mathcal{B}(x_\star, R_1 \beta d^{10} \|x_\star\|)$  and by induction this holds for all  $i \geq T$ . Finally we obtain (A14) by letting  $\rho_1 = (1 - \mu(\tau_1 - \tau_2) \|x_\star\|^2 / 2^{2d})$  and  $\rho_2 = K_3 \tau_3 / (25d^4)$  in (A15).  $\square$

Appendix A.7. Proof of Theorem 3

We begin recalling the following fact on the local Lipschitz property of the generative network  $G$  under the WDC.

**Lemma A10** (Lemma 21 in [16]). Suppose  $x \in \mathcal{B}(x_\star, d\sqrt{\epsilon} \|x_\star\|)$ , and the WDC holds with  $\epsilon < 1 / (200^4 d^6)$ . Then it holds that:

$$\|G(x) - G(x_\star)\| \leq \frac{1.2}{2^d/2} \|x - x_\star\|.$$

We then conclude the proof of Theorem 3 by using the above lemma and the results in the previous sections.

- (I) By assumption  $x_0 \in \mathcal{B}(0, R_\star) \setminus \{0\}$  so that according to Proposition A1 for any  $i \geq 1$  it holds that  $x_i \in \mathcal{B}(0, R_\star) \setminus \{0\}$ .
- (II) By Proposition A3, there exists an integer  $T$  such that  $x_T \in \mathcal{B}^+$  and therefore it satisfies the conclusions of Theorem 3.A

$$\|x_T - x_\star\| \leq C_1 d^{14} \sqrt{\epsilon} \|x_\star\| + C_2 2^d d^{10} \omega \|x_\star\|^{-1}.$$

- (III) Once in  $\mathcal{B}^+$  the iterates of Algorithm 1 converge to  $x_\star$  up to the noise level, as shown by Proposition A4 and Equation (A14)

$$\|x_{i+1} - x_\star\| \leq \rho_1^{i+1-T} \|x_T - x_\star\| + \rho_2 \frac{2^d}{\|x_\star\|} \omega$$

which corresponds to (11) in Theorem 3.B .

- (IV) The reconstruction error (12) in Theorem 3.B, follows then from (11) by applying Lemma A10 and the lower bound (A4).

**Appendix B. Supplementary Proofs**

*Appendix B.1. Supplementary Proofs for Appendix A.2*

Below we prove Lemma A2 on the concentration of the gradient of  $f$  at a differentiable point.

**Proof of Lemma A2.** We begin by noticing that:

$$\bar{v}_x - h_x = [\langle p_x, x \rangle p_x - \langle x, x \rangle \frac{x}{2^{2d}}] + [\langle \tilde{h}_x, x \rangle \frac{\tilde{h}_x}{2^{2d}} - \langle q_x, x \rangle x].$$

Below we show that:

$$\|\langle p_x, x \rangle p_x - \langle x, x \rangle \frac{x}{2^{2d}}\| \leq \frac{50}{2^{2d}} d^3 \sqrt{\epsilon} \max\{\|x\|^2, \|x_\star\|^2\} \|x\|. \tag{A16}$$

and

$$\|\langle q_x, x \rangle p_x - \langle \tilde{h}_x, x \rangle \frac{\tilde{h}_x}{2^{2d}}\| \leq \frac{36}{2^{2d}} d^4 \sqrt{\epsilon} \max\{\|x\|^2, \|x_\star\|^2\} \|x\|. \tag{A17}$$

from which the thesis follows.

Regarding Equation (A16) observe that:

$$\begin{aligned} \|\langle p_x, x \rangle p_x - \langle x, x \rangle \frac{x}{2^{2d}}\| &= \|\langle p_x, x \rangle [p_x - \frac{x}{2^d}] + \langle p_x - \frac{x}{2^d}, \frac{x}{2^d} \rangle x\| \\ &\leq (\| \Lambda_x x \|^2 + \frac{\|x\|^2}{2^d}) \|p_x - \frac{x}{2^d}\| \\ &\leq \frac{50}{2^{2d}} d^3 \sqrt{\epsilon} \|x\|^3 \end{aligned}$$

where in the first inequality we used  $\langle p_x, x \rangle = \| \Lambda_x x \|^2$  and in the second we used Equations (A5) and (A6) of Lemma A1.

Next note that:

$$\begin{aligned} \|\langle q_x, x \rangle q_x - \langle \tilde{h}_x, x \rangle \frac{\tilde{h}_x}{2^{2d}}\| &= \|\langle q_x, x \rangle (q_x - \frac{\tilde{h}_x}{2^d}) + \langle q_x - \frac{\tilde{h}_x}{2^d}, \frac{\tilde{h}_x}{2^d} \rangle \frac{\tilde{h}_x}{2^d}\| \\ &\leq (\|q_x\| + \frac{\|\tilde{h}_x\|}{2^d}) \|x\| \|q_x - \frac{\tilde{h}_x}{2^d}\| \\ &\leq (\frac{13}{12} + 1 + \frac{d}{\pi}) \frac{\|x\| \|x_\star\|}{2^d} \|q_x - \tilde{h}_x\| \\ &\leq \frac{3}{2} d \frac{\|x\| \|x_\star\|}{2^d} \|q_x - \tilde{h}_x\| \end{aligned}$$

where in the second inequality we have the bound (A6) and the definition of  $\tilde{h}_x$ . Equation (A17) is then found by appealing to Equation (A5) in Lemma A1.  $\square$

The previous lemma is now used to control the concentration of the subgradients  $v_x$  of  $f$  around  $h_x$ .

**Proof of Lemma A3.** When  $f$  is differentiable at  $x$ ,  $\nabla f(x) = \bar{v}_x = \bar{v}_x + \eta_x$ , so that by Lemma A2 and the assumption on the noise:

$$\begin{aligned} \|\bar{v}_x - h_x\| &\leq \|\bar{v}_x - h_x\| + \|\eta_x\| \\ &\leq 86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} \max(\|x_\star\|^2, \|x\|^2) \|x\| + \frac{\omega}{2^d} \|x\|. \end{aligned} \tag{A18}$$

Observe, now, that by (9), for any  $x \in \mathbb{R}^k$ ,  $v_x \in \partial f(x) = \text{conv}(v_1, \dots, v_t)$ , and therefore  $v_x = a_1 v_1 + \dots + a_T v_T$  for some  $a_1, \dots, a_T \geq 0$ ,  $\sum_i a_i = 1$ . Moreover for each  $v_i$  there exist a  $w_i$  such that  $v_i = \lim_{\delta \rightarrow 0^+} \bar{v}_{x+\delta w_i}$ . Therefore using Equation (A18), the continuity of  $h_x$  with respect to nonzero  $x$  and  $\sum_i a_i = 1$ :

$$\begin{aligned} \|v_x - h_x\| &\leq \sum_{i=1}^T a_i \|v_i - h_x\| \\ &\leq \sum_{i=1}^T a_i \lim_{\delta \rightarrow 0} \|\tilde{v}_{x+\delta w_i} - h_{x+\delta w_i}\| \\ &\leq 86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} \max(\|x_\star\|^2, \|x\|^2) \|x\| + \frac{\omega}{2^d} \|x\|. \end{aligned}$$

□

The above results are now used to bound the norm of  $v_x \in \partial f(x)$ .

**Proof of Lemma A4.** Since  $\epsilon < 1/(16\pi d^2)^2$  observe that  $86d^4\sqrt{\epsilon} \leq 2d^2$ , therefore by the assumption on the noise level and Lemma A3 it follows that for any  $v_x \in \partial f(x)$  and  $K_2 \leq 1$

$$\|v_x - h_x\| \leq \frac{1}{2^{2d}} \frac{5}{2} d^2 \max(\|x\|^2, \|x_\star\|^2) \|x\|.$$

Next observe that that since  $\|\tilde{h}_x\| \leq d\|x_\star\|$ , we have:

$$\|h_x\| \leq \frac{1}{2^{2d}} \frac{5}{4} d^2 \max(\|x\|^2, \|x_\star\|^2) \|x\|$$

and from  $\|v_x\| \leq \|v_x - h_x\| + \|h_x\|$  we obtain the thesis. □

*Appendix B.2. Supplementary Proofs for Appendix A.3*

In this section we prove Lemma A5 which implies that the norm of the iterates  $x_i$  does not increase in the region  $\mathcal{B}(0, R_\star) \setminus \mathcal{B}(0, C_\star\|x_\star\|)$ .

**Proof of Lemma A5.** Note that the thesis is equivalent to  $2\langle x, v_x \rangle \geq \lambda\mu\|v_x\|^2$ . Next recall that by the WDC for any  $x \in \mathbb{R}^k$  and  $2d\epsilon < 1$ :

$$\frac{(1 - 2\epsilon d)}{2^d} \|x\|^2 \leq \frac{(1 - 2\epsilon)^d}{2^d} \|x\|^2 \leq \|G(x)\|^2 \leq \frac{(1 + 2\epsilon)^d}{2^d} \|x\|^2 \leq \frac{(1 + 4\epsilon d)}{2^d} \|x\|^2.$$

At a nonzero differentiable point  $x \in \mathbb{R}^k \setminus \{0\}$  with  $C_\star\|x_\star\| < \|x\| < R_\star$ , then

$$\begin{aligned} \langle \tilde{v}_x, x \rangle &\geq \|G(x)\|^2 (\|G(x)\|^2 - \|G(x_\star)\|^2) - \|\Lambda_x^\top H \Lambda_x\| \|x\|^2 \\ &\geq \frac{\|x\|^2}{2^{2d}} \left[ (1 - 2\epsilon d)^2 \|x\|^2 - [(1 - 2\epsilon d)(1 + 4\epsilon d) + K_2/d^{4d}] \|x_\star\|^2 \right] \\ &\geq \frac{\|x\|^2}{2^{2d}} \left[ (1 - 4d\epsilon) \|x\|^2 - [(1 + 2d\epsilon) + K_2/d^{4d}] \|x_\star\|^2 \right]. \end{aligned} \tag{A19}$$

Next, by Lemma A4, the definition of the step length  $\mu$  and  $\max(\|x\|^2, \|x_\star\|^2) \leq R_\star^2$ , we have:

$$\frac{\lambda\mu}{2} \|\tilde{v}_x\|^2 \leq \frac{K_3}{2^{2d}} \|x\|^4, \tag{A20}$$

which, using  $\|x\| > C_\star\|x_\star\|$ , gives

$$\begin{aligned} \langle \tilde{v}_x, x \rangle - \frac{\lambda\mu}{2} \|\tilde{v}_x\|^2 &\geq \frac{\|x\|^2}{2^{2d}} \left[ (1 - 4d\epsilon - K_3) \|x\|^2 - (1 + 2d\epsilon + K_2) \|x_\star\|^2 \right] \\ &\geq \frac{\|x\|^2 \|x_\star\|^2}{2^{2d}} \left[ (1 - 4d\epsilon - K_3) C_\star - (1 + 2d\epsilon + K_2) \right]. \end{aligned}$$

We can then conclude by observing that by the assumptions and for small enough constants  $(1 - 4d\epsilon - K_3)C_* - (1 + 2d\epsilon + K_2/d^{44}) > 0$ .

At a non-differentiable point  $x$ , by the characterization of the Clarke subdifferential, we can write  $v_x = \sum_{\ell=1}^m c_\ell v_\ell$  where  $v_\ell = \lim_{\delta_\ell \rightarrow 0^+} \nabla f(x + \delta_\ell w_\ell)$  then

$$\langle x, v_x \rangle = \langle x, \sum_{\ell=1}^m c_\ell v_\ell \rangle = \sum_{\ell=1}^m c_\ell \lim_{\delta_\ell \rightarrow 0^+} \langle x + \delta_\ell w_\ell, \tilde{v}_{x+\delta_\ell w_\ell} \rangle$$

which implies that the lower bound (A19) also holds for  $\langle x, v_x \rangle$ . Similarly Lemma A4 leads to the upper bound (A20) also for  $v_x$ , which then leads to the thesis  $2\langle x, v_x \rangle \geq \lambda\mu\|v_x\|^2$ .  $\square$

Appendix B.3. Supplementary Proofs for Appendix A.4

In the next lemmas we show that  $h_x$  is locally Lipschitz.

**Lemma A11.** For all  $x, y \neq 0$

$$\|h_x - h_y\| \leq \frac{1}{22d} \left[ 2(\|x\|^2 + \|y\|^2) + \left( d^2 + 5d^3 \max\left(\frac{\|x\|}{\|y\|}, \frac{\|y\|}{\|x\|}\right) \right) \|x_*\|^2 \right] \|x - y\|$$

**Proof.** Note that  $\tilde{h}_x = \xi\|x_*\| + \zeta\|x_*\|\hat{x}$  where  $\xi$  and  $\zeta$  are defined in (A23). Then observe that by (A28) and (A29) it follows that  $\|\tilde{h}_x\| \leq d\|x_*\|$ . Furthermore by [16] (Lemma 18) for any nonzero  $x, y \in \mathbb{R}^k$  we have

$$\|\tilde{h}_x - \tilde{h}_y\| \leq \frac{9}{4}d^2 \max\left(\frac{1}{\|x\|}, \frac{1}{\|y\|}\right) \|x_*\| \|x - y\| \tag{A21}$$

Next notice that

$$\|h_x - h_y\| \leq \frac{1}{22d} \|\langle x, x \rangle x - \langle y, y \rangle y\| + \frac{1}{22d} \|\langle \tilde{h}_y, y \rangle \tilde{h}_y - \langle \tilde{h}_x, x \rangle \tilde{h}_x\|,$$

where by triangle inequality the first term on the left hand side can be bounded as

$$\begin{aligned} \|\langle x, x \rangle x - \langle y, y \rangle y\| &\leq (\|x\|^2 \|x - y\| + \left| \|x\|^2 - \|y\|^2 \right| \|y\|) \\ &\leq (\|x\|^2 + \|x\| \|y\| + \|y\|^2) \|x - y\| \\ &\leq 2(\|x\|^2 + \|y\|^2) \|x - y\| \end{aligned}$$

Finally note that by from the bound (A21) we obtain:

$$\begin{aligned} \|\langle \tilde{h}_y, y \rangle \tilde{h}_y - \langle \tilde{h}_x, x \rangle \tilde{h}_x\| &\leq \|\tilde{h}_x\| \|x\| \|\tilde{h}_x - \tilde{h}_y\| + \|\tilde{h}_y\| |\langle \tilde{h}_x, x \rangle - \langle \tilde{h}_y, y \rangle| \\ &\leq \|\tilde{h}_x\| \|x\| \|\tilde{h}_x - \tilde{h}_y\| + \|\tilde{h}_y\| \|\tilde{h}_x\| \|x - y\| + \|\tilde{h}_y\| \|y\| \|\tilde{h}_x - \tilde{h}_y\| \\ &\leq d^2 \|x_*\|^2 \|x - y\| + d(\|x\| + \|y\|) \|x_*\| \|\tilde{h}_x - \tilde{h}_y\| \\ &\leq \left( 2d^2 + 5d^3 \max\left(\frac{\|x\|}{\|y\|}, \frac{\|y\|}{\|x\|}\right) \right) \|x_*\|^2 \|x - y\|. \end{aligned}$$

where we used

$$(\|x\| + \|y\|) \max\left(\frac{1}{\|x\|}, \frac{1}{\|y\|}\right) \leq 2 \frac{\max(\|x\|, \|y\|)}{\min(\|x\|, \|y\|)} = 2 \max\left(\frac{\|x\|}{\|y\|}, \frac{\|y\|}{\|x\|}\right).$$

$\square$

Based on the previous lemma we can now prove that  $h_x$  is locally Lipschitz..

**Lemma A12.** Let  $x \in \mathcal{B}(0, R_*) \setminus \{0\}$ ,  $\lambda \in (0, 1)$ ,  $\mu = K_3 2^{2d} / (8d^4 R_*^2)$  with  $512\pi^2 K_3 < 1$  and  $v_x \in \partial f(x)$ , then for  $x_\lambda = x - \lambda \mu v_x$  it holds that:

$$\|h_x - h_{x_\lambda}\| \leq \frac{7}{16} \|v_x\|$$

**Proof.** Consider  $x \in \mathcal{B}(0, R_*)$  and observe that by Lemma A4,  $d \geq 2$  and the choice of  $\mu$ , for any  $v_x \in \partial f(x)$  and any  $\lambda \in (0, 1)$  we have  $\mu \|v_x\| \leq \|x\|/8$ . It follows that

$$\frac{7}{8} \|x\| \leq \|x_\lambda\| \leq \frac{9}{8} \|x\|, \tag{A22}$$

and in particular

$$\max\left(\frac{\|x\|}{\|x_\lambda\|}, \frac{\|x_\lambda\|}{\|x\|}\right) \leq \frac{8}{7}.$$

Therefore by Lemma A11 we deduce that

$$\begin{aligned} \|h_x - h_{x_\lambda}\| &\leq \frac{\lambda \mu}{2^{2d}} \left[ 2(\|x\|^2 + \|x_\lambda\|^2) + (d^2 + 6d^3) \|x_\star\|^2 \right] \|v_x\| \\ &\leq \frac{\mu}{2^{2d}} (4R_*^2 + (d^2 + 6d^3) \|x_\star\|^2) \|v_x\| \end{aligned}$$

where in the second inequality we used  $\max(\|x\|^2, \|x_\lambda\|^2) \leq R_*^2$  by Proposition A1. The thesis is obtained by substituting the definition of  $\mu$  and using  $K_3 \leq 1$ .  $\square$

Based on the previous result we can now prove Lemma A6.

**Proof of Lemma A6.** Let  $x \in \mathcal{B}(0, R_*) \cap S_\beta^c$ ,

$$\begin{aligned} \|v_x\| &\geq \|h_x\| - \|v_x - h_x\| \\ &\geq \frac{\beta}{2^d} \max(\|x_\star\|^2, \|x\|^2) \|x\| - 86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} \max(\|x_\star\|^2, \|x\|^2) \|x\| - \frac{\omega}{2^d} \|x\| \\ &\geq 6 \frac{\max(\|x\|^2, \|x_\star\|^2) \|x\|}{2^{2d}} (86d^4 \sqrt{\epsilon} + 2^d \omega \|x_\star\|^{-2}) \\ &\geq 6 \left( \frac{\max(\|x\|^2, \|x_\star\|^2) \|x\|}{2^{2d}} 86d^4 \sqrt{\epsilon} + \frac{\omega}{2^d} \|x\| \right) \end{aligned}$$

where we used the definition  $\beta$  in Equation (A7) and Lemma A3.

Next take  $\lambda \in [0, 1]$  and  $x_\lambda = x - \lambda \mu v_x$  with  $v_x \in \partial f(x)$ . Then for any  $v_{x_\lambda} \in \partial f(x_\lambda)$

$$\begin{aligned} \|v_x - v_{x_\lambda}\| &\leq \|v_x - h_x\| + \|h_x - h_{x_\lambda}\| + \|h_{x_\lambda} - v_{x_\lambda}\| \\ &\leq 86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} (\max(\|x_\star\|^2, \|x\|^2) \|x\| + \max(\|x_\star\|^2, \|x_\lambda\|^2) \|x_\lambda\|) \\ &\quad + \frac{\omega}{2^d} (\|x\| + \|x_\lambda\|) + \frac{7}{16} \|v_x\| \\ &\leq 86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} \left(1 + \left(\frac{9}{8}\right)^3\right) \max(\|x_\star\|^2, \|x\|^2) \|x\| + \left(1 + \frac{9}{8}\right) \frac{\omega}{2^d} \|x\| + \frac{7}{16} \|v_x\| \\ &\leq 3 \left(86 \frac{d^4 \sqrt{\epsilon}}{2^{2d}} \max(\|x_\star\|^2, \|x\|^2) \|x\| + \frac{\omega}{2^d} \|x\|\right) + \frac{7}{16} \|v_x\| \\ &\leq \left(\frac{1}{2} + \frac{7}{16}\right) \|v_x\| \end{aligned}$$

where in the first inequality we used Lemma A3 and Lemma A12, in the second inequality (A22) and in the last one (A8).  $\square$

Appendix B.4. Supplementary Proofs for Appendix A.5

Below we prove that the region of  $\mathbb{R}^k$  where we cannot control the norm of the vector field  $h_x$  is contained in two balls around  $x_*$  and  $-\rho_d x_*$ .

We prove Lemma A7 by showing the following.

**Lemma A13.** Suppose  $8\pi d^6 \sqrt{\beta} \leq 1$ . Define:

$$\rho_d := \sum_{i=0}^{d-1} \frac{\sin \check{\theta}_i}{\pi} \left( \prod_{j=i+1}^{d-1} \frac{\pi - \check{\theta}_j}{\pi} \right)$$

where  $\check{\theta}_0 = \pi$  and  $\check{\theta}_i = g(\check{\theta}_{i-1})$ . If  $x \in \mathcal{S}_\beta$ , then we have either:

$$|\bar{\theta}_0| \leq 32d^4 \pi \beta \quad \text{and} \quad \|x\|^2 - \|x_*\|^2 \leq 258\pi \beta d^6 \|x_*\|$$

or

$$|\bar{\theta}_0 - \pi| \leq 8\pi d^4 \sqrt{\beta} \quad \text{and} \quad \|x\|^2 - \rho_d^2 \|x_*\|^2 \leq 281\pi^2 \sqrt{\beta} d^{10} \|x_*\|.$$

In particular, we have:

$$\mathcal{S}_\beta \subset \mathcal{B}(x_*, R_1 \beta d^{10} \|x_*\|) \cup \mathcal{B}(-\rho_d x_*, R_2 \sqrt{\beta} d^{10} \|x_*\|)$$

where  $R_1, R_2$  are numerical constants and  $\rho_d \rightarrow 1$  as  $d \rightarrow \infty$ .

**Proof of Lemma A13.** Without loss of generality, let  $x_* = e_1$  and  $x = r \cos \bar{\theta}_0 \cdot e_1 + r \sin \bar{\theta}_0 \cdot e_2$ , for some  $\bar{\theta}_0 \in [0, \pi]$ , and  $r \geq 0$ . Recall that we call  $\hat{x} = x/\|x\|$  and  $\hat{x}_* = x_*/\|x_*\|$ . We then introduce the following notation:

$$\xi = \prod_{i=0}^{d-1} \frac{\pi - \bar{\theta}_i}{\pi}, \quad \zeta = \sum_{i=0}^{d-1} \frac{\sin \bar{\theta}_i}{\pi} \prod_{j=i+1}^{d-1} \frac{\pi - \bar{\theta}_j}{\pi}, \quad r = \|x\|, \quad R = \max(r^2, 1), \quad (\text{A23})$$

where  $\bar{\theta}_i = g(\bar{\theta}_{i-1})$  with  $g$  as in (A1), and observe that  $\tilde{h}_x = (\xi \hat{x}_* + \zeta \hat{x})$ . Let  $\alpha := \langle \tilde{h}_x, \hat{x} \rangle$ , then we can write:

$$h_x = \frac{1}{2^{2d}} \left[ \langle x, x \rangle x - \langle \tilde{h}_x, x \rangle \tilde{h}_x \right] = \frac{r}{2^{2d}} \left[ r^2 \hat{x} - \alpha (\xi \hat{x}_* + \zeta \hat{x}) \right].$$

Using the definition of  $\hat{x}$  and  $\hat{x}_*$  we obtain:

$$\frac{2^{2d} h_x}{r} = [(r^2 - \alpha \zeta) \cos \bar{\theta}_0 - \alpha \xi] \cdot e_1 + [r^2 - \alpha \zeta] \sin \bar{\theta}_0 \cdot e_2,$$

and conclude that since  $x \in \mathcal{S}_\beta$ , then:

$$|(r^2 - \alpha \zeta) \cos \bar{\theta}_0 - \alpha \xi| \leq \beta R \tag{A24}$$

$$|[r^2 - \alpha \zeta] \sin \bar{\theta}_0| \leq \beta R. \tag{A25}$$

We now list some bounds that will be useful in the subsequent analysis. We have:



$$\bar{\theta}_i \leq \bar{\theta}_{i-1} \text{ for } i \geq 1 \tag{A26}$$

$$\bar{\theta}_i \leq \cos^{-1}(1/\pi) \text{ for } i \geq 2 \tag{A27}$$

$$|\zeta| \leq 1 \tag{A28}$$

$$|\zeta| \leq \frac{d}{\pi} \sin \bar{\theta}_0 \tag{A29}$$

$$\xi \geq \frac{\pi - \bar{\theta}_0}{\pi} d^{-3} \tag{A30}$$

$$\check{\theta}_i \leq \frac{3\pi}{i+3} \text{ for } i \geq 0 \tag{A31}$$

$$\check{\theta}_i \geq \frac{\pi}{i+1} \text{ for } i \geq 0 \tag{A32}$$

$$\bar{\theta}_0 = \pi + O_1(\delta) \Rightarrow |\zeta| \leq \frac{\delta}{\pi} \tag{A33}$$

$$\bar{\theta}_0 = \pi + O_1(\delta) \Rightarrow \zeta = \rho_d + O_1(3d^3\delta) \text{ if } \frac{d^2\delta}{\pi} \leq 1 \tag{A34}$$

$$1/\pi \leq \alpha \leq 1. \tag{A35}$$

The identities (A26) through (A34) can be found in Lemma 16 of [16], while the identity (A35) follows by noticing that  $\alpha = \zeta \cos \bar{\theta}_0 + \xi = \cos \theta_d$  and using (A27) together with  $d \geq 2$ .

**Bound on R.** We now show that if  $x \in \mathcal{S}_\beta$ , then  $r^2 \leq 4d$  and therefore  $R \leq 4d$ .

If  $r^2 \leq 1$ , then the claim is trivial. Take  $r^2 > 1$ , then note that either  $|\sin \bar{\theta}_0| \geq 1/\sqrt{2}$  or  $|\cos \bar{\theta}_0| \geq 1/\sqrt{2}$  must hold. If  $|\sin \bar{\theta}_0| \geq 1/\sqrt{2}$  then from (A25) it follows that  $r^2 - \alpha\zeta \leq \sqrt{2}\beta R = \sqrt{2}\beta r^2$  which implies:

$$r^2 \leq \frac{\alpha\zeta}{1 - \sqrt{2}\beta} \leq \frac{1}{(1 - \sqrt{2}\beta)} \frac{d}{\pi} \leq \frac{d}{2}$$

using (A29) and (A35) in the second inequality and  $\beta < 1/4$  in the third. Next take  $|\cos \bar{\theta}_0| \geq 1/\sqrt{2}$ , then (A24) implies  $|r^2 - \alpha\zeta| \leq \sqrt{2}(\beta r^2 + \alpha\xi)$  which in turn results in:

$$r^2 \leq \frac{\alpha(\zeta + \sqrt{2}\xi)}{1 - \sqrt{2}\beta} \leq 4d$$

using (A28), (A29), (A35) and  $\beta < 1/4$ . In conclusion if  $x \in \mathcal{S}_\beta$  then  $r^2 \leq 4d \Rightarrow R \leq 4d$ .

**Bounds on  $\bar{\theta}_0$ .** We proceed by showing that we only have to analyze the small angle case  $\bar{\theta}_0 \approx 0$  and the large angle case  $\bar{\theta}_0 \approx \pi$ .

At least one of the following three cases must hold:

1.  $\sin \bar{\theta}_0 \leq 16\beta\pi d^4$ : Then we have  $\bar{\theta} = O_1(32\pi\beta\pi d^4)$  or  $\bar{\theta} = \pi + O_1(32\pi\beta\pi d^4)$  as  $32\pi\beta\pi d^4 < 1$ .
2.  $|r^2 - \alpha\zeta| < \sqrt{\beta}R$ : Then (A24), (A35) and  $\beta < 1$  yield  $|\zeta| \leq 2\sqrt{\beta}\pi R$ . Using (A30), we then get  $\bar{\theta} = \pi + O_1(2\sqrt{\beta}\pi^2 d^3 R)$ .
3.  $\sin \bar{\theta}_0 > 16\beta\pi d^4$  and  $|r^2 - \alpha\zeta| \geq \sqrt{\beta}R$ : Then (A25) gives  $|r^2 - \alpha\zeta| \leq \beta M / \sin \bar{\theta}_0$  which used with (A24) leads to:

$$|\alpha\xi| \leq \beta R + |r^2 - \alpha\zeta| \leq \beta R + \frac{\beta R}{\sin \bar{\theta}_0} \leq 2 \frac{\beta R}{\sin \bar{\theta}_0}.$$

Then using (A35), the assumption on  $\sin \bar{\theta}_0$  and  $R \leq 4d$  we obtain  $\xi \leq d^{-3}/2$ . The latter together with (A30) leads to  $\bar{\theta}_0 \geq \pi/2$ . Finally as  $|r^2 - \alpha\zeta| \geq \sqrt{\beta}R$  then (A25)

leads to  $|\sin \bar{\theta}_0| \leq \sqrt{\beta}$ . Therefore as  $\bar{\theta}_0 \geq \pi/2$  and  $\beta < 1$ , we can conclude that  $\bar{\theta}_0 = \pi + O_1(2\sqrt{\beta})$ .

Inspecting the three cases, and recalling that  $R \leq 4d$ , we can see that it suffices to analyze the small angle case  $\bar{\theta}_0 = O_1(32d^4\pi\beta)$  and the large angle case  $\bar{\theta} = \pi + O_1(8\sqrt{\beta}\pi^2d^4)$ .

**Small angle case.** We assume  $\bar{\theta}_0 = O_1(\delta)$  with  $\delta = 32d^4\pi\beta$  and show that  $\|x\|^2 \approx \|x_\star\|^2$ . We begin collecting some bounds. Since  $\bar{\theta}_i \leq \bar{\theta}_0 \leq \delta$ , then  $1 \geq \zeta \geq (1 - \delta/\pi)^d \geq 1 + O_1(2d\delta/\pi)$  assuming  $\delta d/\pi \leq 1/2$ , which holds true since  $64d^5\beta < 1$ . Moreover from (A29) we have  $\zeta = O_1(d\delta/\pi)$ . Finally observe that  $\cos \bar{\theta}_0 = 1 + O_1(\bar{\theta}_0^2/2) = 1 + O_1(\delta/2)$  for  $\delta < 1$ . We then have  $\alpha = 1 + O_1(2d\delta)$  so that  $\alpha\zeta = O_1(d^2\delta)$  and  $\alpha\bar{\zeta} = 1 + O_1(4d^2\delta)$ . We can therefore rewrite (A24) as:

$$(r^2 + O_1(d^2\delta))(1 + O_1(\delta/2)) - (1 + O_1(4d^2\delta)) = O_1(\beta R).$$

Using the bound  $r^2 \leq R \leq 4d$  and the definition of  $\delta$ , we obtain:

$$\begin{aligned} r^2 - 1 &= O_1\left(\frac{\delta r^2}{2} + d^2\delta + \frac{d^2\delta^2}{2} + 4d^2\delta + 4d\beta\right) \\ &= O_1(8d^2\delta + 4d\beta) \\ &= O_1(258\pi d^6\beta) \end{aligned} \tag{A36}$$

**Large angle case.** Here we assume  $\bar{\theta} = \pi + O_1(\delta)$  with  $\delta = 8\sqrt{\beta}\pi^2d^4$  and show that it must be  $\|x\|^2 \approx \rho_d^2\|x_\star\|^2$ .

From (A33) we know that  $\zeta = O_1(\delta/\pi)$ , while from (A34) we know that  $\zeta = \rho_d + O_1(3d^3\delta)$  as long as  $8\sqrt{\beta}\pi d^6 \leq 1$ . Moreover for large angles and  $\delta < 1$ , it holds  $\cos \bar{\theta}_0 = -1 + O_1((\bar{\theta}_0 - \pi)^2/2) = -1 + O_1(\delta^2/2)$ . These bounds lead to:

$$\begin{aligned} \alpha &= \zeta \cos \bar{\theta}_0 + \zeta \\ &= \rho_d + O_1\left(\frac{\delta}{\pi} + \frac{\delta^3}{2\pi} + 3d^3\delta\right) \\ &= \rho_d + O_1(4d^3\delta), \end{aligned}$$

and using  $\rho_d \leq d$ :

$$\begin{aligned} \alpha\zeta &= \rho_d^2 + O_1(4d^3\delta\rho_d + 3d^3\delta\rho_d + 12d^6\delta) = \rho_d^2 + O_1(20d^6\delta), \\ \alpha\bar{\zeta} &= O_1\left(\frac{\delta}{\pi}\rho_d + 4\frac{d^3\delta^2}{\pi}\right) = O_1(2d^3\delta). \end{aligned}$$

Then recall that (A24) is equivalent to  $(r^2 - \alpha\zeta) \cos \bar{\theta}_0 - \alpha\bar{\zeta} = O_1(4\beta d)$ , that is:

$$(r^2 - \rho_d^2 + O_1(20d^6\delta))(1 + O_1(\delta^2/2)) + O_1(2d^3\delta) = O_1(4\beta d)$$

and in particular:

$$\begin{aligned} r^2 - \rho_d^2 &= O_1\left(20d^6\delta + 10d^6\delta^3 + \frac{\rho_d\delta^2}{2} + \frac{r^2\delta^2}{2} + 2d^3\delta + 4\beta d\right) \\ &= O_1(35d^6\delta + 4\beta d) \\ &= O_1(281\pi^2\sqrt{\beta}d^{10}) \end{aligned} \tag{A37}$$

where we used  $\rho_d \leq d$ , the definition of  $\delta$  and  $\delta < 1$ .

**Controlling the distance.** We have shown that it is either  $\bar{\theta}_0 \approx 0$  and  $\|x\|^2 \approx \|x_\star\|^2$  or  $\bar{\theta}_0 \approx \pi$  and  $\|x\|^2 \approx \rho_d^2\|x_\star\|^2$ . We can therefore conclude that it must be either  $x \approx x_\star$  or  $x \approx -\rho_d x_\star$ .

Observe that if a two dimensional point is known to have magnitude within  $\Delta r$  of some  $r$  and is known to be within an angle  $\Delta\theta$  from 0, then its Euclidean distance to the point of coordinates  $(r, 0)$  is no more that  $\Delta r + (r + \Delta r)\Delta\theta$ . Similarly we can write:

$$\|x - x_\star\| \leq \| \|x\| - \|x_\star\| \| + (\|x_\star\| + \| \|x\| - \|x_\star\| \|) \bar{\theta}_0. \tag{A38}$$

In the small angle case, by (A36), (A38), and  $\|x_\star\| \| \|x\| - \|x_\star\| \| \leq \| \|x\|^2 - \|x_\star\|^2 \|$ , we have:

$$\|x - x_\star\| \leq 258\pi d^6 \beta + (1 + 258\pi d^6 \beta) 32d^4 \pi \beta \leq 550 \pi d^{10} \beta.$$

Next we notice that  $\rho_2 = 1/\pi$  and  $\rho_d \geq \rho_{d-1}$  as follows from the definition and (A31), (A32). Then considering the large angle case and using (A37) we have:

$$\| \|x\| - \rho_d \| \leq \frac{281\pi^2 \sqrt{\beta} d^{10}}{\|x\| + \rho_d} \leq 281\pi^3 \sqrt{\beta} d^{10}.$$

The latter, together with (A38), yields:

$$\begin{aligned} \|x + \rho_d x_\star\| &\leq \| \|x\| - \rho_d \| + (\rho_d + \| \|x\| - \rho_d \|) (\pi - \bar{\theta}_0) \\ &\leq 281\pi^3 \sqrt{\beta} d^{10} + (d + 281\pi^3 \sqrt{\beta} d^{10}) 8\sqrt{\beta} \pi^2 d^4 \\ &\leq 284\pi^3 \sqrt{\beta} d^{10} \end{aligned}$$

where in the second inequality we have used  $\rho_d \leq d$  and in the third  $8\sqrt{\beta} \pi^2 d^4 \leq 1$ .

We conclude by noticing that  $\rho_d \rightarrow 1$  as  $d \rightarrow 1$  as shown in [16] (Lemma 16).  $\square$

We next will show that the values of the loss function in a neighborhood of  $x_\star$  are strictly smaller than those in a neighborhood of  $-\rho_d x_\star$ .

Recall that  $f(x) := 1/4 \|G(x)G(x)^\top - G(x_\star)G(x_\star)^\top - H\|_F^2$ , we next define the following loss functions:

$$\begin{aligned} f_0(x) &:= \frac{1}{4} \|G(x)G(x)^\top - G(x_\star)G(x_\star)^\top\|_F^2, \\ f_H(x) &:= f_0(x) - \frac{1}{2} \langle G(x)G(x)^\top - G(x_\star)G(x_\star)^\top, H \rangle_F, \\ f_E(x) &:= \frac{1}{2^{2d+2}} \left( \|x\|^4 + \|x_\star\|^4 - 2\langle x, \tilde{h}_x \rangle^2 \right). \end{aligned}$$

In particular notice that  $f(x) = f_H(x) + 1/4 \|H\|_F^2$ . Below we show that assuming the WDC is satisfied,  $f_0(x)$  concentrates around  $f_E(x)$ .

**Lemma A14.** *Suppose that  $d \geq 2$  and the WDC holds with  $\epsilon < 1/(16\pi d^2)^2$ , then for all nonzero  $x, x_\star \in \mathbb{R}^k$*

$$|f_0(x) - f_E(x)| \leq \frac{16}{2^{2d}} (\|x\|^4 + \|x_\star\|^4) d^4 \sqrt{\epsilon}$$

**Proof.** Observe that:

$$\begin{aligned} |f_0(x) - f_E(x)| &\leq \frac{1}{4} \left| \|G(x)\|^4 - \frac{1}{2^{2d}} \|x\|^4 \right| \\ &\quad + \frac{1}{4} \left| \|G(x_\star)\|^4 - \frac{1}{2^{2d}} \|x_\star\|^4 \right| \\ &\quad + \frac{1}{2} \left| \langle G(x), G(x_\star) \rangle^2 - \frac{1}{2^{2d}} \langle x, \tilde{h}_x \rangle^2 \right|. \end{aligned}$$

We analyze each term separately. The first term can be bounded as:

$$\begin{aligned} \frac{1}{4} \left| \|G(x)\|^4 - \frac{1}{2^{2d}} \|x\|^4 \right| &= \frac{1}{4} \left| \|G(x)\|^2 + \frac{1}{2^d} \|x\|^2 \right| \left| \|G(x)\|^2 - \frac{1}{2^d} \|x\|^2 \right| \\ &\leq \frac{1}{4} \frac{1}{2^d} \left( \frac{13}{12} + 1 \right) \|x\|^2 \left| \|G(x)\|^2 - \frac{1}{2^d} \|x\|^2 \right| \\ &\leq \frac{1}{4} \frac{1}{2^d} \left( \frac{13}{12} + 1 \right) \|x\|^2 \cdot 24 \frac{d^3 \sqrt{\epsilon}}{2^d} \|x\|^2 \\ &\leq \frac{1}{2^{2d}} 13d^3 \sqrt{\epsilon} \|x\|^4 \end{aligned}$$

where in the first inequality we used (A6) and in the second inequality (A5). Similarly we can bound the second term:

$$\frac{1}{4} \left| \|G(x_*)\|^4 - \frac{1}{2^{2d}} \|x_*\|^4 \right| \leq \frac{1}{2^{2d}} 13d^3 \sqrt{\epsilon} \|x_*\|^4.$$

We next note that  $\|\tilde{h}_x\| \leq (1 + d/\pi)\|x_*\|$  and therefore from (A6) and  $d \geq 2$  we obtain:

$$\left| \|G(x)\| \|G(x_*)\| + \|x\| \frac{\|\tilde{h}_x\|}{2^d} \right| \leq \frac{1}{2^d} \left( \frac{13}{12} + 1 + \frac{d}{\pi} \right) \|x\| \|x_*\| \leq \frac{1}{2^d} \frac{3}{2} d \|x\| \|x_*\|$$

We can then conclude that:

$$\begin{aligned} \frac{1}{2} \left| \langle G(x), G(x_*) \rangle^2 - \langle x, \frac{\tilde{h}_x}{2^d} \rangle^2 \right| &\leq \frac{1}{2} \left| \langle x, \Lambda_x^T \Lambda_{x_*} x_* - \tilde{h}_x \rangle \right| \left| \|G(x)\| \|G(x_*)\| + \|x\| \frac{\|\tilde{h}_x\|}{2^d} \right| \\ &\leq \frac{1}{2} \|x\| 24 \frac{d^3 \sqrt{\epsilon}}{2^d} \|x_*\| \frac{1}{2^d} \frac{3}{2} d \|x\| \|x_*\| \\ &\leq \frac{9}{2^{2d}} d^4 \sqrt{\epsilon} (\|x_*\|^4 + \|x\|^4) \end{aligned}$$

□

We next consider the loss  $f_E$  and show that in a neighborhood  $-\rho_d x_*$ , this loss function has larger values than in a neighborhood of  $x_*$ .

**Lemma A15.** Fix  $0 < a \leq 1/(2\pi^3 d^3)$  and  $\phi_d \in [\rho_d, 1]$  then:

$$\begin{aligned} f_E(x) &\leq \frac{1}{2^{2d+2}} \|x_*\|^4 + \frac{1}{2^{2d+2}} [(a + \phi_d)^4 - 2\phi_d^2 + 2\pi da] \|x_*\|^4 \quad \forall x \in \mathcal{B}(\phi_d x_*, a \|x_*\|) \text{ and} \\ f_E(x) &\geq \frac{1}{2^{2d+2}} \|x_*\|^4 + \frac{1}{2^{2d+2}} [(a - \phi_d)^4 - 2\rho_d^2 \phi_d^2 - 40\pi d^3 a] \|x_*\|^4 \quad \forall x \in \mathcal{B}(-\phi_d x_*, a \|x_*\|). \end{aligned}$$

**Proof.** Let  $x \in \mathcal{B}(\phi_d x_*, a \|x_*\|)$  then observe that  $0 \leq \bar{\theta}_i \leq \bar{\theta}_0 \leq \pi a / 2\phi_d$  and  $(\phi_d - a)\|x_*\| \leq \|x\| \leq (a + \phi_d)\|x_*\|$ . Then observe that:

$$\begin{aligned} \langle x, \frac{\tilde{h}_d}{2^d} \rangle &= \frac{1}{2^d} \left( \prod_{i=0}^{d-1} \frac{\pi - \bar{\theta}_i}{\pi} \right) \|x_*\| \|x\| \cos \bar{\theta}_0 + \frac{1}{2^d} \sum_{i=0}^{d-1} \frac{\sin \bar{\theta}_i}{\pi} \prod_{j=i+1}^{d-1} \frac{\pi - \bar{\theta}_j}{\pi} \|x_*\| \|x\| \\ &\geq \frac{1}{2^d} \left( \prod_{i=0}^{d-1} \frac{\pi - \frac{\pi a}{2\phi_d}}{\pi} \right) (\phi_d - a) \|x_*\|^2 \left( 1 - \frac{\pi^2 a^2}{8\phi_d^2} \right) \\ &\geq \frac{1}{2^d} \left( 1 - \frac{da}{\phi_d} \right) (\phi_d - a) \left( 1 - \frac{\pi^2 a^2}{8\phi_d^2} \right) \|x_*\|^2. \end{aligned}$$

using  $\cos \theta \geq 1 - \theta^2/2$  and  $(1 - x)^d \geq (1 - 2dx)$  as long as  $0 \leq x \leq 1$ . We can therefore write:

$$\begin{aligned}
 f_E(x) - \frac{\|x_\star\|^4}{2^{2d+2}} &\leq \frac{1}{2^{2d+2}} \|x\|^4 - \frac{1}{2^{2d+1}} \left(1 - \frac{da}{\phi_d}\right)^2 (\phi_d - a)^2 \left(1 - \frac{\pi^2 a^2}{8\phi_d^2}\right)^2 \|x_\star\|^4 \\
 &\leq \frac{1}{2^{2d+2}} \left[ (\phi_d + a)^4 - 2\left(1 - 2\frac{da}{\phi_d}\right) (\phi_d - a)^2 \left(1 - \frac{\pi^2 a^2}{4\phi_d^2}\right) \right] \|x_\star\|^4
 \end{aligned}$$

where in the second inequality we used  $(1 - x)^2 \geq 1 - 2x$  for all  $x \in \mathbb{R}$ . We then observe that:

$$\begin{aligned}
 \left(1 - 2\frac{da}{\phi_d}\right) (\phi_d - a)^2 \left(1 - \frac{\pi^2 a^2}{4\phi_d^2}\right) &\geq \left(1 - \frac{\pi^2 a^2}{4\phi_d^2} - \frac{2ad}{\phi_d}\right) \phi_d^2 + a(a - 2\phi_d) \left(1 - 2\frac{da}{\phi_d}\right) \left(1 - \frac{\pi^2 a^2}{4\phi_d^2}\right) \\
 &\geq \phi_d^2 - a\left(\frac{1}{2\pi d^3} + 2d\phi_d\right) + a(a - 2\phi_d) \left(1 - 2\frac{da}{\phi_d}\right) \left(1 - \frac{\pi^2 a^2}{4\phi_d^2}\right) \\
 &\geq \phi_d^2 - a\left(\frac{1}{2\pi d^3} + 2d\phi_d + 2\phi_d\right) \\
 &\geq \phi_d^2 - \pi da,
 \end{aligned}$$

where in the second inequality we have used  $\pi^3 d^3 a \leq 2$  and in the last one  $d \geq 2$  and  $\phi_d \leq 1$ . We can then conclude that:

$$f_E(x) - \frac{\|x_\star\|^4}{2^{2d+2}} \leq \frac{1}{2^{2d+2}} [(\phi_d + a)^4 - 2(\phi_d^2 - \pi da)] \|x_\star\|^4$$

We next take  $x \in \mathcal{B}(-\phi_d x_\star, a \|x_\star\|)$  which implies  $0 \leq \pi - \bar{\theta}_0 \leq \pi^2 a/2 =: \delta$  and  $\|x\| \leq (a + \phi_d) \|x_\star\|$ . We then note that for  $\tilde{\zeta}$  and  $\zeta$  as defined in (A23) we have:

$$\begin{aligned}
 |x^\top \tilde{h}_x|^2 &\leq (|\tilde{\zeta}| + |\zeta|)^2 (a + \phi_d)^2 \|x_\star\|^4 \\
 &\leq \left(\frac{\delta}{\pi} + 3d^3\delta + \rho_d\right)^2 (a + \phi_d)^2 \|x_\star\|^4 \\
 &\leq \left(\frac{\pi^3 d^3}{2} a + \rho_d\right)^2 (a + \phi_d)^2 \|x_\star\|^4 \\
 &\leq (2\pi^3 d^3 a + \rho_d^2) (a + \phi_d)^2 \|x_\star\|^4 \\
 &\leq 20\pi d^3 a + \rho_d^2 \phi_d^2
 \end{aligned}$$

where the second inequality is due to (A33) and (A34), the rest from  $d \geq 2$ ,  $\rho_d \leq \phi_d \leq 1$  and  $2\pi^3 d^3 a \leq 1$ . Finally using  $(\phi_d - a) \|x_\star\| \leq \|x\|$ , we can then conclude that:

$$f_E(x) - \frac{\|x_\star\|^4}{2^{2d+2}} \geq \frac{1}{2^{2d+2}} [(\phi_d - a)^4 - 2(20\pi d^3 a + \rho_d^2 \phi_d^2)] \|x_\star\|^4.$$

□

The above two lemmas are now used to prove Lemma A8.

**Proof of Proposition A8.** Let  $x \in \mathcal{B}(\pm\phi_d x_\star, \varphi \|x_\star\|)$  for a  $0 < \varphi < 1$  that will be specified below, and observe that by the assumptions on the noise:

$$\begin{aligned}
 |(G(x)G(x)^\top - G(x_\star)G(x_\star)^\top), H)_F| &\leq |G(x)^\top H G(x)| + |G(x_\star)^\top H G(x_\star)| \\
 &\leq \frac{\omega}{2^d} (\|x\|^2 + \|x_\star\|^2) \\
 &\leq \frac{\omega}{2^d} ((\phi_d + \varphi)^2 + 1) \|x_\star\|^2,
 \end{aligned}$$

and therefore by Lemma A14:

$$\begin{aligned}
 |f_0(x) - f_E(x)| + \frac{1}{2} |\langle G(x)G(x)^\top - G(x_*)G(x_*)^\top, H \rangle_F| &\leq \\
 &\leq \frac{16}{2^{2d}} ((\phi_d + \varphi)^4 + 1) \|x_*\|^4 d^4 \sqrt{\epsilon} + \frac{\omega}{2^d} ((\phi_d + \varphi)^2 + 1) \|x_*\|^2 \\
 &\leq \frac{272}{2^{2d}} \|x_*\|^4 d^4 \sqrt{\epsilon} + \frac{\omega}{2^d} ((\phi_d + \varphi)^2 + 1) \|x_*\|^2
 \end{aligned}$$

We next take  $\varphi = \epsilon$  and  $x \in \mathcal{B}(\phi_d x_*, \varphi \|x_*\|)$ , so that by Lemma A15 and the assumption  $2^d d^{44} \omega \leq K_2 \|x_*\|^2$ , we have:

$$\begin{aligned}
 f_H(x) &\leq f_E(x) + |f_0(x) - f_E(x)| + \frac{1}{2} |\langle G(x)G(x)^\top - G(x_*)G(x_*)^\top, H \rangle_F| \\
 &\leq \frac{1}{2^{2d+2}} [1 + (\epsilon + \phi_d)^4 - 2\phi_d^2 + 2\pi d \epsilon] \|x_*\|^4 + 272d^4 \sqrt{\epsilon} \|x_*\|^4 + \frac{\omega}{2^{d+1}} (2 + 2\epsilon + \epsilon^2) \|x_*\|^2 \\
 &\leq \frac{1}{2^{2d+2}} [1 - 2\phi_d^2 + (\epsilon + \phi_d)^4] \|x_*\|^4 + \frac{1}{2^{2d}} \left( \frac{3}{2} 2^d \|x_*\|^{-2} \omega + \frac{\pi d}{2} + 272d^4 \right) \sqrt{\epsilon} \|x_*\|^4 + \frac{\omega}{2^d} \|x_*\|^2 \\
 &\leq \frac{1}{2^{2d+2}} [1 - 2\phi_d^2 + (\epsilon + \phi_d)^4] \|x_*\|^4 + \frac{1}{2^{2d}} \left( \frac{3}{2} K_2 d^{-44} + \frac{\pi d}{2} + 272d^4 \right) \sqrt{\epsilon} \|x_*\|^4 + K_2 \frac{\|x_*\|^4}{2^{2d}} d^{-44}.
 \end{aligned}$$

Similarly if  $y \in \mathcal{B}(-\phi_d x_*, \varphi \|x_*\|)$ , and  $\varphi = \epsilon$  we obtain:

$$\begin{aligned}
 f_H(y) &\geq f_E(y) - |f_0(y) - f_E(y)| - \frac{1}{2} |\langle G(y)G(y)^\top - G(x_*)G(x_*)^\top, H \rangle_F| \\
 &\geq \frac{1}{2^{2d+2}} [1 - 2\phi_d^2 \rho_d^2 + (\epsilon - \phi_d)^4] \|x_*\|^4 - \frac{1}{2^{2d}} \left( \frac{3}{2} 2^d \|x_*\|^{-2} \omega + 10\pi d^3 + 272d^4 \right) \sqrt{\epsilon} \|x_*\|^4 - \frac{\omega}{2^d} \|x_*\|^2 \\
 &\geq \frac{1}{2^{2d+2}} [1 - 2\phi_d^2 \rho_d^2 + (\epsilon - \phi_d)^4] \|x_*\|^4 - \frac{1}{2^{2d}} \left( \frac{3}{2} K_2 d^{-44} + 10\pi d^3 + 272d^4 \right) \sqrt{\epsilon} \|x_*\|^4 - K_2 \frac{\|x_*\|^4}{2^{2d}} d^{-44}.
 \end{aligned}$$

In order to guarantee that  $f(y) > f(x)$ , it suffices to have:

$$2(1 - \rho_d^2)\phi_d^2 - 8K_2 d^{-44} > 4C_d \sqrt{\epsilon}$$

with  $C_d := (544d^4 + 10\pi d^3 \pi + 3K_2 d^{-44} + \pi d/2 + 1/100)$ , that is to require:

$$\varphi = \epsilon < \left( \frac{(1 - \rho_d^2)\phi_d^2 - 4K_2 d^{-44}}{2C_d} \right)^2.$$

Finally notice that by Lemma 17 in [16] it holds that  $1 - \rho_d \geq (K(d + 2))^{-2}$  for some numerical constant  $K$ , we therefore choose  $\varphi = \varrho/d^{12}$  for some  $\varrho > 0$  small enough.  $\square$

### Appendix B.5. Supplementary Proofs for Appendix A.6

In this section we use strong convexity and smoothness to prove convergence to  $x_*$  up to the noise variance  $\omega$ . The idea is to show that every vector in the subgradient points in the direction  $(x - x_*)$ . Recall that the gradient in the noiseless case was:

$$\bar{v}_x = \Lambda_x^\top [\Lambda_x x x^\top \Lambda_x^\top - \Lambda_{x_*} x_* x_*^\top \Lambda_{x_*}^\top] \Lambda_x x.$$

We show that by continuity of  $\Lambda_x$ , when  $x$  is close to  $x_*$ , then  $\bar{v}_x$  it is close to:

$$\bar{\bar{v}}_x := \Lambda_x^\top \Lambda_x [x x^\top - x_* x_*^\top] \Lambda_x^\top \Lambda_x x.$$

which in turn concentrates around:

$$\check{v}_x := \frac{1}{2^d} [x x^\top - x_* x_*^\top] x$$

by the WDC.

We begin by recalling the following result which can be found in the proof of Lemma 22 of [16].



**Lemma A16.** Suppose  $x \in \mathcal{B}(x_*, d\sqrt{\epsilon}\|x_*\|)$  and the WDC holds with  $\epsilon < 1/(200^4d^6)$ . Then it holds that:

$$\|\Lambda_x^\top \Lambda_x x_* - \Lambda_{x_*}^\top \Lambda_{x_*} x_*\| \leq \frac{1}{16} \frac{1}{2^d} \|x - x_*\|.$$

We now prove that  $\bar{v}_x \approx \bar{\bar{v}}_x$  for  $x$  close to  $x_*$ .

**Lemma A17.** Suppose  $x \in \mathcal{B}(x_*, d\sqrt{\epsilon}\|x_*\|)$  and the WDC holds with  $\epsilon < 1/(200^4d^6)$ . Then it holds that:

$$\|\bar{v}_x - \bar{\bar{v}}_x\| \leq \frac{13}{96} \frac{1}{2^{2d}} \|x\| \|x_*\| \|x - x_*\|$$

**Proof.** Let  $g_{x,*} := \Lambda_x^\top \Lambda_x x_*$  and  $q_{x,*} := \Lambda_{x_*}^\top \Lambda_{x_*} x_*$ . Then observe that:

$$\begin{aligned} \|\bar{v}_x - \bar{\bar{v}}_x\| &= \|\langle g_{x,*}, x \rangle g_{x,*} - \langle q_{x,*}, x \rangle q_{x,*}\| \\ &\leq (\|x\| \|g_{x,*}\| + \|\Lambda_x x\| \|\Lambda_{x_*} x_*\|) \|g_{x,*} - q_{x,*}\| \\ &\leq \frac{13}{6} \frac{1}{2^d} \|x\| \|x_*\| \|g_{x,*} - q_{x,*}\| \\ &\leq \frac{13}{96} \frac{1}{2^{2d}} \|x\| \|x_*\| \|x - x_*\| \end{aligned}$$

where the second inequality follows from Lemma A1, and the third from Lemma A16.  $\square$

We next prove that by the WDC  $\bar{\bar{v}}_x \approx \check{v}_x$ .

**Lemma A18.** Suppose the WDC holds with  $\epsilon < 1/(16\pi d^2)^2$ . Then for all nonzero  $x, x_*$ :

$$\|\bar{\bar{v}}_x - \check{v}_x\| \leq \frac{25}{12} \frac{4\epsilon d}{2^{2d}} \|x\| \|x - x_*\| \|x + x_*\|$$

**Proof.** For notational convenience we define  $E_d := I_k/2^d$  the scaled identity in  $\mathbb{R}^d$ ,  $Q_x := \Lambda_x^\top \Lambda_x$  and  $M_{x,*} = xx^\top - x_*x_*^\top$ . Next observe that:

$$\begin{aligned} \|\bar{\bar{v}}_x - \check{v}_x\| &\leq \|(Q_x - E_d)M_{x,*}Q_x x\| + \|E_d M_{x,*}(Q_x - E_d)x\|, \\ &\leq \frac{25}{12} \frac{1}{2^d} \|Q_x - E_d\| \|M_{x,*}\| \|x\| \\ &\leq \frac{25}{12} \frac{4\epsilon d}{2^{2d}} \|M_{x,*}\| \|x\| \end{aligned}$$

where the second inequality follows from Lemma A1 and the third from (17) in [15]. We conclude by noticing that  $\|M_{x,*}\| \leq \|x - x_*\| \|x + x_*\|$ .  $\square$

Next consider the quartic function  $\check{f}(x) := 1/2^{2d+2} \|xx^\top - x_*x_*^\top\|_F^2$  and observe that:

$$\begin{aligned} \nabla \check{f}(x) &= \frac{1}{2^{2d}} (xx^\top - x_*x_*^\top)x = \check{v}_x, \\ \nabla^2 \check{f}(x) &= \frac{1}{2^{2d}} (\|x\|^2 I_n + 2xx^\top - x_*x_*^\top). \end{aligned}$$

Following [63] we show that  $\check{v}_x$  is  $\beta$ -smooth and strongly convex, in turn deriving the result in Lemma A19.

**Lemma A19.** Assume  $\|x - x_*\| \leq \gamma \|x_*\|$  with  $\gamma < 1/5$ . Take  $\tau \geq 3(1 + \gamma)^2 + 1$ , then:

$$\|\check{v}_x - \frac{\tau}{2^{2d}} \|x_*\|^2 (x - x_*)\| \leq \sqrt{\tau(\tau - \alpha)} \|x - x_*\| \frac{\|x_*\|^2}{2^{2d}}$$

where  $\alpha = 2 - 9\gamma$ .

**Proof.** Let  $\|x - x_\star\| \leq \gamma\|x_\star\|$  with  $\gamma < 1/5$ , then:

$$(1 - \gamma)\|x_\star\| \leq \|x\| \leq (1 + \gamma)\|x_\star\| \tag{A39}$$

$$\|\Delta\|\|x\| \leq \gamma(1 + \gamma)\|x_\star\|^2 \tag{A40}$$

since  $\Delta = x - x_\star$  satisfies  $\|\Delta\| \leq \|x\|$ . Using (A40) we then obtain

$$\begin{aligned} xx^T &= x_\star x_\star^T + \Delta x^T + x \Delta^T - \Delta \Delta^T, \\ &\succeq x_\star x_\star^T - 3\|\Delta\|\|x\|I_k, \\ &\succeq x_\star x_\star^T - 3\gamma(1 + \gamma)\|x_\star\|^2 I_k. \end{aligned}$$

We therefore have:

$$\begin{aligned} 2^{2d}\nabla^2\check{f}(x) &= \|x\|^2 I_k + 2xx^T - x_\star x_\star^T \\ &\succeq \|x\|^2 I_k + x_\star x_\star^T - 6\gamma(1 + \gamma)\|x_\star\|^2 I_k, \\ &\succeq ((1 - \gamma)^2 + 1 - 6\gamma(1 + \gamma))\|x_\star\|^2 I_k, \\ &\succeq (2 - 9\gamma)\|x_\star\|^2 I_k, \end{aligned}$$

where in the second line we have used (A39) and in the third  $0 < \gamma < 1/5$ .

Next notice that using (A39) we have

$$2^{2d}\|\nabla^2\check{f}(x)\| \leq 3\|x\|^2 + \|x_\star\|^2 \leq (3(1 + \gamma)^2 + 1)\|x_\star\|^2$$

and therefore, for all  $x \in \mathcal{B}(x_\star, \gamma\|x_\star\|)$

$$(2 - 9\gamma)\|x_\star\|^2 I_k \preceq 2^{2d}\nabla^2\check{f}(x) \preceq (3(1 + \gamma)^2 + 1)\|x_\star\|^2.$$

The above bounds imply in particular that for all  $x \in \mathcal{B}(x_\star, \gamma\|x_\star\|^2)$  the gradient of  $\check{f}$  satisfies the *regularity condition* (see [63]):

$$2\langle \nabla\check{f}(x), x - x_\star \rangle \geq \mu\|\nabla\check{f}(x)\|^2 + \lambda\|x - x_\star\|^2, \tag{A41}$$

where  $\lambda = (2 - 9\gamma)\|x_\star\|^2/2^{2d}$  and  $1/\mu = (3(1 + \gamma)^2 + 1)\|x_\star\|^2/2^{2d}$ . By (A41) we can then conclude that for  $\sigma \geq 1/\mu$ :

$$\begin{aligned} \|\nabla\check{f}(x) - \sigma(x - x_\star)\|^2 &\leq \|\nabla\check{f}(x)\|^2 + \sigma^2\|x - x_\star\|^2 - 2\sigma\langle \nabla\check{f}(x), x - x_\star \rangle \\ &\leq (1 - \sigma\mu)\|\nabla\check{f}(x)\|^2 + \sigma(\sigma - \lambda)\|x - x_\star\|^2 \\ &\leq \sigma(\sigma - \lambda)\|x - x_\star\|^2 \end{aligned}$$

Finally letting  $\sigma = \tau\|x_\star\|^2/2^d$ ,  $\alpha = (2 - 9\gamma)$ , and by the assumptions on  $\gamma$  and  $\tau$  the thesis follows.  $\square$

We finally can prove Lemma A9

**Proof of Lemma A9.** Let  $x \in \mathcal{B}(x_\star, d\sqrt{\epsilon}\|x_\star\|)$ , then  $\|x + x_\star\| \leq (2 + d\sqrt{\epsilon})\|x_\star\|$  and since  $\epsilon < 1/(200^4 d^6)$  we have from Lemma A17 and Lemma A18

$$\|\bar{v}_x - \check{v}_x\| \leq \|\bar{v}_x - \bar{v}_x\| + \|\bar{v}_x - \check{v}_x\| \leq \frac{1}{2^{2d}} \frac{13}{48} \|x\| \|x_\star\| \|x - x_\star\|.$$

If  $x \in \mathcal{B}(x_\star, d\sqrt{\epsilon}\|x_\star\|)$  is a differentiable point of  $f$ , then by Lemma A19 and the assumption (13) on the noise  $H$ , it holds that:

$$\|\bar{v}_x - \frac{\tau}{2^{2d}}\|x_\star\|^2(x - x_\star)\| \leq \left[\frac{13}{48}(1 + \gamma) + \sqrt{\tau(\tau - \alpha)}\right] \frac{\|x_\star\|^2}{2^{2d}} \|x - x_\star\| + \frac{\omega}{2^d}(1 + \gamma)\|x_\star\|.$$

where  $\gamma = d\sqrt{\epsilon}$ ,  $\alpha = 2 - 9\gamma$  and  $\tau \geq 3(1 + \gamma)^2 + 1$ .

In general, if  $x \in \mathcal{B}(x_*, \gamma\|x_*\|)$  and  $v_x \in \partial f(x)$ , by the characterization of Clarke subdifferential (9) and the previous results:

$$\begin{aligned} \|v_x - \frac{\tau}{2^{2d}}\|x_*\|^2(x - x_*)\| &\leq \sum_{t=1}^T c_t \|v_t - \frac{\tau}{2^{2d}}\|x_*\|^2(x - x_*)\| \\ &\leq \left[ \frac{13}{48}(1 + \gamma) + \sqrt{\tau(\tau - \alpha)} \right] \frac{\|x_*\|^2}{2^{2d}} \|x - x_*\| + \frac{\omega}{2^d}(1 + \gamma)\|x_*\|. \end{aligned}$$

We finally obtain the thesis by noting that by the assumptions on  $\epsilon$  it holds that  $\gamma = d\sqrt{\epsilon} < 1/(400)^2$  and taking  $\tau = \tau_1 = 21/5$ ,  $\tau_2 = 17/5$  and  $\tau_3 = (1 + 1/(400)^2)$ .  $\square$

### Appendix C. Proof of Theorem 1

By Theorem 3 it suffices to show that with high probability the weight matrices  $\{W_i\}_{i=1}^d$  satisfy the WDC, and that  $\omega/2^d$  upper bounds the spectral norm of the noise term  $\Lambda_x^\top H \Lambda_x$  where

- in the **Spiked Wishart Model**  $H = \Sigma_N - \Sigma$  where  $\Sigma_N = Y^\top Y/N$  and the  $\Sigma = y_* y_*^\top + \sigma^2 I_n$ ;
- in the **Spiked Wigner Model**  $H = v\mathcal{H}$  where  $\mathcal{H} \sim GOE(n)$ .

Regarding the Weight Distribution Condition (WDC), we observe that it was initially proposed in [15], where it was shown to hold with high probability for networks with random Gaussian weights under an expansivity condition on their dimensions. It was later shown in [62] that a less restrictive expansivity rate is sufficient.

**Lemma A20** (Theorem 3.2 in [62]). *There are constants  $C, c > 0$  with the following property. Let  $0 < \epsilon < 1$  and suppose  $W \in \mathbb{R}^{n \times k}$  has i.i.d.  $\mathcal{N}(0, 1/n)$  entries. Suppose that  $n \geq ck\epsilon^{-2} \log(1/\epsilon)$ . Then with probability at least  $1 - \exp(-Ck)$ ,  $W$  satisfies the WDC with constant  $\epsilon$ .*

By a union bound over all layers, using the above result we can conclude that the WDC holds simultaneously for all layers of the network with probability at least  $1 - \sum_{i=1}^d e^{-Cn_i-1}$ . Note in particular that this argument does not requires the independence of the weight matrices  $\{W_i\}_{i=1}^d$ .

By Lemma A20, with high probability the random generative network  $G$  satisfies the WDC. Therefore if we can guarantee that the assumptions on the noise term  $H$  are satisfied, then the proof of the main Theorem 1 follows from the deterministic Theorem 3 and Lemma A20.

Before turning to the bounds of the noise terms in the spiked models, we recall the following lemma which bounds the number of possible  $\Lambda_x$  for  $x \neq 0$ . Note that this is related to the number of possible regions defined by a deep ReLU network.

**Lemma A21.** *Consider a network  $G$  as defined in (3) with  $d \geq 2$ , weight matrices  $W_i \in \mathbb{R}^{n_i \times n_{i-1}}$  with i.i.d. entries  $\mathcal{N}(0, 1/n_i)$ . Then, with probability one, for any  $x \neq 0$  the number of different matrices  $\Lambda_x$  is*

$$|\{\Lambda_x | x \neq 0\}| \leq 10^{d^2} (n_1^d n_2^{d-1} \dots n_d)^k \leq (n_1^d n_2^{d-1} \dots n_d)^{8k}$$

**Proof.** The first inequality follows from Lemma 16 and the proof of Lemma 17 in [15]. For the second inequality notice that as  $k \geq 1$ ,  $n_1 > k$  and  $d \geq 2$  it follows that  $7k \log(n_1^d n_2^{d-1} \dots n_d) \geq 7k d(d + 1) \log(n_1)/2 \geq 3.5d(d + 1) \log(2) \geq d^2 \log(10)$ .  $\square$

In the next section we use this lemma to control the noise term  $\Lambda_x^\top H \Lambda_x$  on the event where the WDC holds.

Appendix C.1. Spiked Wigner Model

Recall that in the Wigner model  $Y = G(x_*)G(x_*)^T + v\mathcal{H}$  and the symmetric noise matrix  $\mathcal{H}$  follows a *Gaussian Orthogonal Ensemble*  $\text{GOE}(n)$ , that is  $\mathcal{H}_{ii} \sim \mathcal{N}(0, 2/n)$  for all  $1 \leq i \leq n$  and  $\mathcal{H}_{ji} = \mathcal{H}_{ij} \sim \mathcal{N}(0, 1/n)$  for  $1 \leq j < i \leq n$ . Our goal is to bound  $\|\Lambda_x^T \mathcal{H} \Lambda_x\|$  uniformly over  $x$  with high probability.

Fix  $x \in \mathbb{R}^k$ , and let  $\mathcal{N}_{1/4}$  be a 1/4-net on the sphere  $S^{k-1}$  such that (see for example [64])  $|\mathcal{N}_{1/4}| \leq 9^k$  and

$$\|\Lambda_x^T \mathcal{H} \Lambda_x\| \leq 2 \max_{z \in \mathcal{N}_{1/4}} |\langle \Lambda_x^T \mathcal{H} \Lambda_x z, z \rangle|.$$

Observe next that for any  $v \in \mathbb{R}^n$  by the definition of  $\text{GOE}(n)$  it holds that

$$v^T H v = \sum_i^n H_{ii} v_i^2 + 2 \sum_{i < j}^n H_{ij} v_i v_j \sim \mathcal{N}\left(0, 2\left(\sum_i v_i^4 + 2 \sum_{i < j} v_i^2 v_j^2\right)/n\right) = \mathcal{N}\left(0, 2\|v\|^4/n\right).$$

Therefore for any  $z \in \mathcal{N}_{1/4}$  let  $\ell_{x,z} := \Lambda_x z \in \mathbb{R}^n$ , then  $\ell_{x,z}^T \mathcal{H} \ell_{x,z} \sim \mathcal{N}(0, 2\|\ell_{x,z}\|^4/n)$ . In particular by (A6), the quadratic form  $\ell_{x,z}^T \mathcal{H} \ell_{x,z}$  is sub-Gaussian with parameter  $\gamma^2$  given by

$$\gamma^2 := \frac{2}{n} \left(\frac{13}{12}\right)^2 \frac{1}{2^{2d}}.$$

Then for fixed  $x \in \mathbb{R}^k$ , standard sub-Gaussian tail bounds (e.g., [64]) and a union bound over  $\mathcal{N}_{1/4}$  give for any  $u \geq 0$

$$\begin{aligned} \mathbb{P}[\|\Lambda_x^T \mathcal{H} \Lambda_x\| \geq 2u] &\leq \mathbb{P}\left[\max_{z \in \mathcal{N}_{1/4}} \|\ell_{x,z}^T \mathcal{H} \ell_{x,z}\| \geq u\right] \\ &\leq \sum_{z \in \mathcal{N}_{1/4}} \mathbb{P}[\|\ell_{x,z}^T \mathcal{H} \ell_{x,z}\| \geq u] \leq 2 \cdot 9^k e^{-\frac{u^2}{2\gamma^2}}. \end{aligned}$$

Lemma A21, then ensures that the number of possible  $\Lambda_x$  is at most  $(n_1^d n_2^{d-1} \dots n_d)^{8k}$ , so a union bound over this set allows us to conclude that

$$\begin{aligned} \mathbb{P}[\|\Lambda_x^T \mathcal{H} \Lambda_x\| \leq 2u, \text{ for all } x] &\geq 1 - (n_1^d n_2^{d-1} \dots n_d)^{8k} \mathbb{P}[\|\Lambda_x^T \mathcal{H} \Lambda_x\| \geq 2u] \\ &\geq 1 - 2 \exp(8k \log(2 n_1^d n_2^{d-1} \dots n_d) - u^2/(2\gamma^2)) \end{aligned}$$

Choosing then  $u = \sqrt{2\gamma^2 \cdot 9k \log(2 n_1^d n_2^{d-1} \dots n_d)}$  and substituting the definition of  $\gamma^2$ , we obtain

$$\mathbb{P}[\|\Lambda_x^T \mathcal{H} \Lambda_x\| \leq \frac{1}{2^d} \sqrt{\frac{169k \log(2 n_1^d n_2^{d-1} \dots n_d)}{n}}, \text{ for all } x] \geq 1 - 2e^{-k \log(2 n_1^d n_2^{d-1} \dots n_d)}$$

which implies the thesis as  $n = n_d$  and  $\log(n) \leq \log(2 n_1^d n_2^{d-1} \dots n_d)$ .

Appendix C.2. Spiked Wishart Model

Each row  $\{y_i\}_{i=1}^N$  of the matrix  $Y$  in (1) can be seen as i.i.d. samples from  $\mathcal{N}(0, \Sigma)$  where  $\Sigma = y_* y_*^T + \sigma^2 I_n$ . In the minimization problem (4) we take  $M = \Sigma_N - \sigma^2 I_n$  where  $\Sigma_N$  is the empirical covariance matrix  $Y^T Y/N$ . The symmetric noise matrix  $H$  is then given by  $H = \Sigma_N - \Sigma$  and by the Law of Large Numbers  $H \rightarrow 0$  as  $N \rightarrow \infty$ . We bound  $\|\Lambda_x^T H \Lambda_x\|$  with high probability uniformly over  $x \in \mathbb{R}^k$ .

Fix  $x \in \mathbb{R}^k$ , let  $\mathcal{N}_{1/4}$  be a 1/4-net on the sphere  $S^{k-1}$  such that  $|\mathcal{N}_{1/4}| \leq 9^k$ , and notice that:

$$\|\Lambda_x^T H \Lambda_x\| \leq 2 \max_{z \in \mathcal{N}_{1/4}} |z^T \Lambda_x^T H \Lambda_x z|.$$

By a union bound on  $\mathcal{N}_{1/4}$  we obtain for any fixed  $z \in \mathcal{N}_{1/4}$ :

$$\mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \geq 2u] \leq 9^k \mathbb{P}[|z^T \Lambda_x^T H \Lambda_x z| \geq u].$$

Let  $\ell_{x,z} := \Lambda_x z$  and  $s_i := \ell_{x,z}^T y_i$ , so that we can write

$$z^T \Lambda_x^T H \Lambda_x z = \frac{1}{N} \sum_{i=1}^N ((\ell_{x,z}^T y_i)^2 - \mathbb{E}[(\ell_{x,z}^T y_i)^2]) = \frac{1}{N} \sum_{i=1}^N (s_i^2 - \mathbb{E}[s_i^2])$$

and in particular

$$\mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \geq 2u] \leq 9^k \mathbb{P}[|z^T \Lambda_x^T H \Lambda_x z| \geq u] = 9^k \mathbb{P}\left[\frac{1}{N} \left| \sum_{i=1}^N s_i^2 - \mathbb{E}[s_i^2] \right| \geq u\right]$$

Observe then that  $s_i \sim \mathcal{N}(0, \gamma^2)$  with  $\gamma^2 = \ell_{x,z}^T \Sigma \ell_{x,z}$ . It follows for  $u \in [0, \gamma^2]$  by the small deviation bound for  $\chi^2$  random variables (e.g., [33] (Example 2.11))

$$\mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \geq 2u] \leq 2 \exp\left(2k \log 3 - \frac{Nu^2}{8\gamma^4}\right)$$

Recall now that  $|\{\Lambda_x | x \neq 0\}| \leq (n_1^d n_2^{d-1} \dots n_d)^{8k}$ , then proceeding as for the Wigner case by a union bound over all possible  $\Lambda_x$ :

$$\begin{aligned} \mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \leq 2u, \text{ for all } x] &\geq 1 - (n_1^d n_2^{d-1} \dots n_d)^{8k} \mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \geq 2u] \\ &\geq 1 - 2 \exp\left(8k \log(2 n_1^d n_2^{d-1} \dots n_d) - \frac{Nu^2}{8\gamma^4}\right). \end{aligned}$$

Substituting  $u = \sqrt{8\gamma^4 \cdot 9k \log(2 n_1^d n_2^{d-1} \dots n_d) / N}$  we find that:

$$\mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \leq 2\sqrt{\frac{72k \log(2 n_1^d n_2^{d-1} \dots n_d)}{N}} \gamma^2, \text{ for all } x] \geq 1 - 2e^{-k \log(n)} \tag{A42}$$

since  $\log(n) \leq \log(2 n_1^d n_2^{d-1} \dots n_d)$ .

Similarly if  $u > \gamma^2$  by large deviation bounds for sub-exponential variables

$$\mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \leq 2u, \text{ for all } x] \geq 1 - 2 \exp\left(8k \log(2 n_1^d n_2^{d-1} \dots n_d) - \frac{Nu}{8\gamma^2}\right).$$

Substituting  $u = 8\gamma^2 \cdot 9k \log(2 n_1^d n_2^{d-1} \dots n_d) / N$  we find that:

$$\mathbb{P}[\|\Lambda_x^T H \Lambda_x\| \leq 2\frac{72k \log(3 n_1^d n_2^{d-1} \dots n_d)}{N} \gamma^2, \text{ for all } x] \geq 1 - 2e^{-k \log(n)} \tag{A43}$$

Finally observe that using (A6) for bounding  $\|\Lambda_x\|^2$  (by the WDC) and  $\sigma^2 + \|y_\star\|^2$  for bounding  $\|\Sigma\|$ , we have

$$\gamma^2 \leq \frac{13}{12} \frac{1}{2^d} (\|y_\star\|^2 + \sigma^2),$$

which combined with (A42) and (A43) implies the thesis.

## References

1. Johnstone, I.M. On the distribution of the largest eigenvalue in principal components analysis. *Ann. Stat.* **2001**, *29*, 295–327. [[CrossRef](#)]
2. Amini, A.A.; Wainwright, M.J. High-dimensional analysis of semidefinite relaxations for sparse principal components. In Proceedings of the 2008 IEEE International Symposium on Information Theory, Toronto, ON, Canada, 6–11 July 2008; pp. 2454–2458.
3. Deshpande, Y.; Montanari, A. Sparse PCA via covariance thresholding. In Proceedings of the Advances in Neural Information Processing Systems, Montréal, QC, Canada, 8–13 December 2014; pp. 334–342.
4. Vu, V.; Lei, J. Minimax rates of estimation for sparse PCA in high dimensions. In Proceedings of the 15th International Conference on Artificial Intelligence and Statistics, La Palma, Canary Islands, Spain, 21–23 April 2012; pp. 1278–1286.
5. Abbe, E.; Bandeira, A.S.; Bracher, A.; Singer, A. Decoding binary node labels from censored edge measurements: Phase transition and efficient recovery. *IEEE Trans. Netw. Sci. Eng.* **2014**, *1*, 10–22. [[CrossRef](#)]
6. Bandeira, A.S.; Chen, Y.; Lederman, R.R.; Singer, A. Non-unique games over compact groups and orientation estimation in cryo-EM. *Inverse Probl.* **2020**, *36*, 064002. [[CrossRef](#)]
7. Javanmard, A.; Montanari, A.; Ricci-Tersenghi, F. Phase transitions in semidefinite relaxations. *Proc. Natl. Acad. Sci. USA* **2016**, *113*, E2218–E2223. [[CrossRef](#)] [[PubMed](#)]
8. McSherry, F. Spectral partitioning of random graphs. In Proceedings of the 42nd IEEE Symposium on Foundations of Computer Science, Newport Beach, CA, USA, 8–11 October 2001; pp. 529–537.
9. Deshpande, Y.; Abbe, E.; Montanari, A. Asymptotic mutual information for the binary stochastic block model. In Proceedings of the 2016 IEEE International Symposium on Information Theory (ISIT), Barcelona, Spain, 10–15 July 2016; pp. 185–189.
10. Moore, C. The computer science and physics of community detection: Landscapes, phase transitions, and hardness. *arXiv* **2017**, arXiv:1702.00467.
11. D’Aspremont, A.; Ghaoui, L.; Jordan, M.; Lanckriet, G. A direct formulation for sparse PCA using semidefinite programming. *Adv. Neural Inf. Process. Syst.* **2004**, *17*, 41–48. [[CrossRef](#)]
12. Berthet, Q.; Rigollet, P. Optimal detection of sparse principal components in high dimension. *Ann. Stat.* **2013**, *41*, 1780–1815. [[CrossRef](#)]
13. Bandeira, A.S.; Perry, A.; Wein, A.S. Notes on computational-to-statistical gaps: predictions using statistical physics. *arXiv* **2018**, arXiv:1803.11132.
14. Kunisky, D.; Wein, A.S.; Bandeira, A.S. Notes on computational hardness of hypothesis testing: Predictions using the low-degree likelihood ratio. *arXiv* **2019**, arXiv:1907.11636.
15. Hand, P.; Voroninski, V. Global guarantees for enforcing deep generative priors by empirical risk. *IEEE Trans. Inf. Theory* **2019**, *66*, 401–418. [[CrossRef](#)]
16. Heckel, R.; Huang, W.; Hand, P.; Voroninski, V. Rate-optimal denoising with deep neural networks. *arXiv* **2018**, arXiv:1805.08855.
17. Hand, P.; Leong, O.; Voroninski, V. Phase retrieval under a generative prior. In Proceedings of the Advances in Neural Information Processing Systems, Montréal, QC, Canada, 3–8 December 2018; pp. 9136–9146.
18. Ma, F.; Ayaz, U.; Karaman, S. Invertibility of convolutional generative networks from partial measurements. *Adv. Neural Inf. Process. Syst.* **2018**, *31*, 9628–9637.
19. Hand, P.; Joshi, B. Global Guarantees for Blind Demodulation with Generative Priors. In Proceedings of the Advances in Neural Information Processing Systems, Vancouver, BC, Canada, 8–14 December 2019; pp. 11535–11543.
20. Song, G.; Fan, Z.; Lafferty, J. Surfing: Iterative optimization over incrementally trained deep networks. In Proceedings of the Advances in Neural Information Processing Systems, Vancouver, BC, Canada, 8–14 December 2019; pp. 15034–15043.
21. Bora, A.; Jalal, A.; Price, E.; Dimakis, A.G. Compressed sensing using generative models. In Proceedings of the 34th International Conference on Machine Learning, Sydney, Australia, 6–11 August 2017; Volume 70, pp. 537–546.
22. Asim, M.; Shamsad, F.; Ahmed, A. Blind Image Deconvolution using Deep Generative Priors. *arXiv* **2019**, arXiv:cs.CV/1802.04073.
23. Hand, P.; Leong, O.; Voroninski, V. Compressive Phase Retrieval: Optimal Sample Complexity with Deep Generative Priors. *arXiv* **2020**, arXiv:2008.10579.
24. Hand, P.; Voroninski, V. Compressed sensing from phaseless gaussian measurements via linear programming in the natural parameter space. *arXiv* **2016**, arXiv:1611.05985.
25. Li, X.; Voroninski, V. Sparse signal recovery from quadratic measurements via convex programming. *SIAM J. Math. Anal.* **2013**, *45*, 3019–3033. [[CrossRef](#)]
26. Ohlsson, H.; Yang, A.Y.; Dong, R.; Sastry, S.S. Compressive phase retrieval from squared output measurements via semidefinite programming. *arXiv* **2011**, arXiv:1111.6323.
27. Cai, T.; Li, X.; Ma, Z. Optimal rates of convergence for noisy sparse phase retrieval via thresholded Wirtinger flow. *Ann. Stat.* **2016**, *44*, 2221–2251. [[CrossRef](#)]
28. Wang, G.; Zhang, L.; Giannakis, G.B.; Akçakaya, M.; Chen, J. Sparse phase retrieval via truncated amplitude flow. *IEEE Trans. Signal Process.* **2017**, *66*, 479–491. [[CrossRef](#)]
29. Yuan, Z.; Wang, H.; Wang, Q. Phase retrieval via sparse wirtinger flow. *J. Comput. Appl. Math.* **2019**, *355*, 162–173. [[CrossRef](#)]
30. Aubin, B.; Loureiro, B.; Maillard, A.; Krzakala, F.; Zdeborová, L. The spiked matrix model with generative priors. In Proceedings of the Advances in Neural Information Processing Systems, Vancouver, BC, Canada, 8–14 December 2019; pp. 8366–8377.



31. Cocola, J.; Hand, P.; Voroninski, V. Nonasymptotic Guarantees for Spiked Matrix Recovery with Generative Priors. In Proceedings of the Advances in Neural Information Processing Systems, Vancouver, BC, Canada, 11 December 2020; Volume 33.
32. Johnstone, I.M.; Lu, A.Y. On consistency and sparsity for principal components analysis in high dimensions. *J. Am. Stat. Assoc.* **2009**, *104*, 682–693. [[CrossRef](#)] [[PubMed](#)]
33. Wainwright, M.J. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*; Cambridge University Press: Cambridge, UK, 2019; Volume 48.
34. Montanari, A.; Richard, E. Non-negative principal component analysis: Message passing algorithms and sharp asymptotics. *IEEE Trans. Inf. Theory* **2015**, *62*, 1458–1484. [[CrossRef](#)]
35. Deshpande, Y.; Montanari, A.; Richard, E. Cone-constrained principal component analysis. *Adv. Neural Inf. Process. Syst.* **2014**, *27*, 2717–2725.
36. Zou, H.; Hastie, T.; Tibshirani, R. Sparse principal component analysis. *J. Comput. Graph. Stat.* **2006**, *15*, 265–286. [[CrossRef](#)]
37. Krauthgamer, R.; Nadler, B.; Vilenchik, D.; others. Do semidefinite relaxations solve sparse PCA up to the information limit? *Ann. Stat.* **2015**, *43*, 1300–1322. [[CrossRef](#)]
38. Berthet, Q.; Rigollet, P. Computational lower bounds for Sparse PCA. *arXiv* **2013**, arXiv:1304.0828.
39. Cai, T.; Ma, Z.; Wu, Y. Sparse PCA: Optimal rates and adaptive estimation. *Ann. Stat.* **2013**, *41*, 3074–3110. [[CrossRef](#)]
40. Ma, T.; Wigderson, A. Sum-of-squares lower bounds for sparse PCA. *Adv. Neural Inf. Process. Syst.* **2015**, *28*, 1612–1620.
41. Lesieur, T.; Krzakala, F.; Zdeborová, L. Phase transitions in sparse PCA. In Proceedings of the 2015 IEEE International Symposium on Information Theory (ISIT), Hong Kong, China, 14–19 June 2015; pp. 1635–1639.
42. Brennan, M.; Bresler, G. Optimal average-case reductions to sparse pca: From weak assumptions to strong hardness. *arXiv* **2019**, arXiv:1902.07380.
43. Arous, G.B.; Wein, A.S.; Zadik, I. Free energy wells and overlap gap property in sparse PCA. In Proceedings of the Conference on Learning Theory, PMLR, Graz, Austria, 9–12 July 2020; pp. 479–482.
44. Fan, J.; Liu, H.; Wang, Z.; Yang, Z. Curse of heterogeneity: Computational barriers in sparse mixture models and phase retrieval. *arXiv* **2018**, arXiv:1808.06996.
45. Richard, E.; Montanari, A. A statistical model for tensor PCA. *Adv. Neural Inf. Process. Syst.* **2014**, *27*, 2897–2905.
46. Decelle, A.; Krzakala, F.; Moore, C.; Zdeborová, L. Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications. *Phys. Rev. E* **2011**, *84*, 066106. [[CrossRef](#)] [[PubMed](#)]
47. Perry, A.; Wein, A.S.; Bandeira, A.S.; Moitra, A. Message-Passing Algorithms for Synchronization Problems over Compact Groups. *Commun. Pure Appl. Math.* **2018**, *71*, 2275–2322. [[CrossRef](#)]
48. Oymak, S.; Jalali, A.; Fazel, M.; Eldar, Y.C.; Hassibi, B. Simultaneously structured models with application to sparse and low-rank matrices. *IEEE Trans. Inf. Theory* **2015**, *61*, 2886–2908. [[CrossRef](#)]
49. Dhar, M.; Grover, A.; Ermon, S. Modeling sparse deviations for compressed sensing using generative models. *arXiv* **2018**, arXiv:1807.01442.
50. Shah, V.; Hegde, C. Solving linear inverse problems using gan priors: An algorithm with provable guarantees. In Proceedings of the 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Calgary, AB, Canada, 15–20 April 2018; pp. 4609–4613.
51. Mixon, D.G.; Villar, S. Sunlayer: Stable denoising with generative networks. *arXiv* **2018**, arXiv:1803.09319.
52. Yeh, R.A.; Chen, C.; Lim, T.Y.; Schwing, A.G.; Hasegawa-Johnson, M.; Do, M.N. Semantic image inpainting with deep generative models. *arXiv* **2016**, arXiv:1607.07539.
53. Sønderby, C.K.; Caballero, J.; Theis, L.; Shi, W.; Huszár, F. Amortised map inference for image super-resolution. *arXiv* **2016**, arXiv:1610.04490.
54. Yang, G.; Yu, S.; Dong, H.; Slabaugh, G.; Dragotti, P.L.; Ye, X.; Liu, F.; Arridge, S.; Keegan, J.; Guo, Y.; et al. DAGAN: Deep de-aliasing generative adversarial networks for fast compressed sensing MRI reconstruction. *IEEE Trans. Med Imaging* **2017**, *37*, 1310–1321. [[CrossRef](#)]
55. Qiu, S.; Wei, X.; Yang, Z. Robust One-Bit Recovery via ReLU Generative Networks: Improved Statistical Rates and Global Landscape Analysis. *arXiv* **2019**, arXiv:1908.05368.
56. Xue, Y.; Xu, T.; Zhang, H.; Long, L.R.; Huang, X. Segan: Adversarial network with multi-scale l1 loss for medical image segmentation. *Neuroinformatics* **2018**, *16*, 383–392. [[CrossRef](#)] [[PubMed](#)]
57. Heckel, R.; Hand, P. Deep Decoder: Concise Image Representations from Untrained Non-convolutional Networks. In Proceedings of the International Conference on Learning Representations, New Orleans, LA, USA, 6–9 May 2019.
58. Heckel, R.; Soltanolkotabi, M. Denoising and regularization via exploiting the structural bias of convolutional generators. *arXiv* **2019**, arXiv:1910.14634.
59. Heckel, R.; Soltanolkotabi, M. Compressive sensing with un-trained neural networks: Gradient descent finds the smoothest approximation. *arXiv* **2020**, arXiv:2005.03991.
60. Aubin, B.; Loureiro, B.; Baker, A.; Krzakala, F.; Zdeborová, L. Exact asymptotics for phase retrieval and compressed sensing with random generative priors. In Proceedings of the First Mathematical and Scientific Machine Learning Conference, PMLR, Princeton, NJ, USA, 20–24 July 2020; pp. 55–73.
61. Clason, C. Nonsmooth Analysis and Optimization. *arXiv* **2017**, arXiv:1708.04180.

- 
62. Daskalakis, C.; Rohatgi, D.; Zampetakis, M. Constant-Expansion Suffices for Compressed Sensing with Generative Priors. *arXiv* **2020**, arXiv:2006.04237.
  63. Chi, Y.; Lu, Y.M.; Chen, Y. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Trans. Signal Process.* **2019**, *67*, 5239–5269. [[CrossRef](#)]
  64. Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*; Cambridge University Press: Cambridge, UK, 2018; Volume 47.