
Can Neural Networks Achieve Optimal Computational-statistical Tradeoff? An Analysis on Single-Index Model

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we tackle the following question: Can neural networks trained with
2 gradient-based methods achieve the optimal statistical-computational tradeoff
3 in learning Gaussian single-index models? Prior research has shown that any
4 polynomial-time algorithm under the statistical query (SQ) framework requires
5 $\Omega(d^{s^*/2} \vee d)$ samples, where s^* is the generative exponent representing the intrinsic
6 difficulty of learning the underlying model. However, it remains unknown whether
7 neural networks can achieve this sample complexity. Inspired by prior techniques
8 such as label transformation and landscape smoothing for learning single-index
9 models, we propose a unified gradient-based algorithm for training a two-layer
10 neural network in polynomial time. Our method is adaptable to a variety of loss and
11 activation functions, covering a broad class of existing approaches. We show that
12 our algorithm learns a feature representation that strongly aligns with the unknown
13 signal θ^* , with sample complexity $\tilde{O}(d^{s^*/2} \vee d)$, matching the SQ lower bound up
14 to a polylogarithmic factor for all generative exponents $s^* \geq 1$. Furthermore, we
15 extend our approach to the setting where θ^* is k -sparse for $k = o(\sqrt{d})$ by introduc-
16 ing a novel weight perturbation technique that leverages the sparsity structure. We
17 derive a corresponding SQ lower bound of order $\tilde{\Omega}(k^{s^*})$, matched by our method
18 up to a polylogarithmic factor. Our framework, especially the weight perturba-
19 tion technique, is of independent interest, and suggests potential gradient-based
20 solutions to other problems such as sparse tensor PCA.

21 1 Introduction

22 The success of neural networks is largely attributed to their remarkable ability to learn rich and
23 useful features from data during gradient-based training (Girshick et al., 2014). This feature-learning
24 capability allows them to outperform traditional methods like kernel-based approaches, which rely on
25 predefined features (Allen-Zhu and Li, 2019; Ghorbani et al., 2019; Refinetti et al., 2021). However,
26 when trained using (stochastic) gradient descent, neural networks can sometimes fall into a “kernel
27 regime”, where their behavior resembles that of a fixed kernel method, constrained by their random
28 initialization (Jacot et al., 2018; Chizat et al., 2019). In this regime, the ability of the network to learn
29 complex representations is severely limited, undermining the primary advantage of deep learning.
30 Therefore, it is crucial to understand when and how neural networks trained with gradient-based
31 method can perform effective feature learning to unlock their full potential, particularly in scenarios
32 where a balance between computational efficiency and statistical performance is essential.

33 In this work, we approach this question in the context of Gaussian single-index models, a canonical
34 class of problems in statistics and learning (MacCullagh and Nelder, 1989; Ichimura, 1993; Hristache

et al., 2001; Härdle et al., 2004). The model is defined as follows: for covariates $z \sim \mathcal{N}(0, I_d)$, the output y depends on the inner product $\langle \theta^*, z \rangle$ with an unknown signal $\theta^* \in \mathbb{R}^d$ through a link distribution p , i.e., $y \sim p(\cdot | \langle \theta^*, z \rangle)$. The goal is to recover θ^* using i.i.d. samples $(z_1, y_1), \dots, (z_n, y_n)$ generated by the underlying model. While $n = \Omega(d)$ samples suffice to recover θ^* information-theoretically (Bach, 2017; Damian et al., 2024), achieving this efficiently is difficult for polynomial-time algorithms, where the required sample size also depends on properties of the link distribution p , creating a computational-statistical gap. For example, when y is a polynomial of $\langle \theta^*, z \rangle$, it has been shown that two-layer neural networks with square loss need $d^{\Theta(q^*)}$ samples (Arous et al., 2021; Bietti et al., 2022; Damian et al., 2023), where q^* is the information exponent of the polynomial link function (Arous et al., 2021; Dudeja and Hsu, 2018). Such sample complexity is indeed inevitable under the correlational statistical query (CSQ) framework, leading to a computational-statistical gap for $q^* \geq 2$.

However, the CSQ framework does not capture the fundamental limits of all gradient-based algorithms. Recent works have shown that by leveraging higher-order terms in the gradient, neural networks can learn polynomials with as few as $\tilde{O}(d)$ samples (Lee et al., 2024; Arnaboldi et al., 2024). It turns out that the intrinsic learning difficulty is captured by another quantity called the *generative exponent* s^* , which is at most 2 for polynomial link functions, and the corresponding SQ lower bound on the sample complexity is $n = \Omega(d^{s^*/2})$ (Damian et al., 2024).¹ Thus, there is no computational-statistical gap for learning polynomial single-index models using neural networks. However, for general single-index models with $s^* \geq 3$, no gradient-based algorithm for neural networks has been shown to match the SQ lower bound, leaving it an open problem (Arnaboldi et al., 2024; Lee et al., 2024).

Furthermore, learning the Gaussian single-index model can benefit from additional structures in the signal θ^* , such as sparsity, which can significantly reduce the sample complexity compared to those depending on the ambient dimension d (Candès et al., 2006; Donoho et al., 2009; Raskutti et al., 2012). For example, in matrix PCA, the best rank-1 estimator achieves a near-optimal sample complexity of $\tilde{O}(d)$ due to the BBP transition (Baik et al., 2005; Choo and d’Orsi, 2021). However, under extreme sparsity, sparse estimators require $\tilde{O}(k^2)$ samples using methods like diagonal thresholding (Johnstone and Lu, 2009) or semidefinite relaxation (d’Aspremont et al., 2004), which improves upon the $\tilde{O}(d)$ sample complexity but exhibits a unique computational-statistical gap from the information-theoretic lower bound $\Omega(k \log d)$ (Wang et al., 2016). For sparse single-index models with information exponent $q^* = 1$, gradient descent on diagonal linear networks nearly achieves the information-theoretic lower bound thanks to its implicit regularization effect (Fan et al., 2023). Nonetheless, how to achieve the optimal sample complexity for general $s^* \geq 1$ is also unknown under the sparse setting.

Contributions. Towards characterizing the fundamental feature learning capability of neural networks in the Gaussian single-index model, our main contributions are as follows:

1. We propose a unified recipe of gradient-based algorithms for training a two-layer neural network to learn the Gaussian single-index model. Our method integrates a general gradient oracle with a weight perturbation technique, carefully designed to exploit the underlying structure of the Gaussian single-index model. This allows the neural network to perform feature learning of the unknown signal θ^* in a computationally efficient manner. Our framework encompasses many existing approaches as special cases, such as batch reusing (Dandi et al., 2024; Lee et al., 2024), label transformation (Chen and Meka, 2020), and landscape smoothing (Damian et al., 2023).
2. We show that for an *unknown link distribution* p with any generative exponent $s^* \geq 1$, the weights of the neural network achieve strong recovery of the true signal θ^* after training by our algorithm using $\tilde{O}(d^{s^*/2} \vee d)$ samples and polynomial running time. Our method achieves the SQ lower bound up to a polylogarithmic factor, and is the first gradient-based algorithm for training two-layer neural networks that attains the nearly optimal computational-statistical tradeoff for Gaussian single-index models with any $s^* \geq 1$. Figure 1 (a) illustrates an example for $s^* = 4$.
3. Furthermore, our method is able to take advantage of additional structural information of the true signal θ^* . Specifically, we consider the case where θ^* is k -sparse for $k = o(\sqrt{d})$, and develop a *novel weight perturbation procedure* tailored to the sparsity of θ^* . Equipped with this, we

¹This $\Omega(d^{s^*/2})$ sample complexity lower bound is essentially for the detection problem. Dudeja and Hsu (2021) shows that there is an estimation-detection gap for tensor PCA under the SQ framework, though it is unclear whether such gap exists universally. Throughout the paper, we always refer to the SQ lower bound as the detection lower bound, since detection in general is assumed to be easier than estimation.

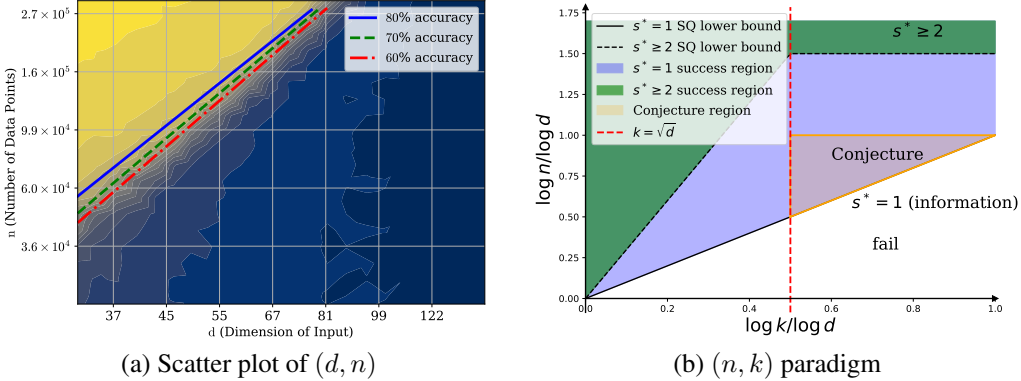


Figure 1: (a) The contour plots of $(\log d, \log n, \text{acc}(d, n))$ for Algorithm 1 under model $y = \langle z, \theta^* \rangle^2 \exp(-\langle z, \theta^* \rangle^2)$, which has generative exponent $s^* = 4$ (Example 2.2). Here $\text{acc}(d, n)$ is the average of the largest 8 values of the alignment between the neuron weights and the unknown signal θ^* . The slopes of these contour lines are all close to 2, indicating a sample complexity $n \approx d^2$ for $s^* = 4$. (b) The paradigm of sample complexity achieved by our algorithm for different generative exponent s^* and sparsity level k , illustrating the success of achieving computational-statistical tradeoff.

86 show that the weights of the neural network can achieve strong recovery of the sparse signal θ^*
 87 after training with $\tilde{O}(k^{s^*})$ samples in polynomial time for any generative exponent $s^* \geq 1$. This
 88 sample complexity is also nearly optimal according to the sample complexity lower bound we
 89 establish for SQ algorithms, which might be of independent interest. Also, our method suggests a
 90 new approach to achieve the computational-statistical tradeoff for sparse tensor PCA.

91 In summary, our work provides a unified framework for training neural networks that can achieve
 92 the nearly optimal computational-statistical tradeoff for the Gaussian single-index model with any
 93 generative exponent $s^* \geq 1$. Our method not only tackles the intrinsic difficulty of learning the under-
 94 lying model posed by the link distribution p , but also leverages the additional structural information
 95 of the true signal θ^* that benefits the learning process. Integrating these results, our method attains
 96 nearly optimal balance between computational efficiency and statistical performance across almost all
 97 regimes of sparsity levels and generative exponent $s^* \geq 1$, as illustrated in Figure 1 (b).

98 **Related Works.** Our work contributes to the recent research on the computation-statistical tradeoff
 99 in learning single-index models. The information-theoretic limit for estimating the latent signal is
 100 $n = \Omega(d)$ (Bach, 2017; Damian et al., 2024), but the sample complexity lower bound varies across
 101 computational models, potentially revealing a computational-statistical gap.

102 The information exponent q^* (Dudeja and Hsu, 2018; Arous et al., 2021) governs the sample com-
 103 plexity for learning Gaussian single-index models in the CSQ framework (Chen et al., 2020; Bietti
 104 et al., 2022; Damian et al., 2022; Dandi et al., 2023; Abbe et al., 2023; Ba et al., 2023). Notably,
 105 Arous et al. (2021) show that online SGD has a sample complexity of $n = \tilde{O}(d^{q^*-1})$, which is worse
 106 than the CSQ lower bound $n = \Omega(d^{q^*/2})$ (Abbe et al., 2023; Damian et al., 2022). This gap can be
 107 closed by a loss landscape smoothing technique (Damian et al., 2023) originally developed for tensor
 108 PCA (Anandkumar et al., 2017; Biroli et al., 2020). Our work extends beyond the CSQ framework,
 109 aligning with more general SQ algorithms (Feldman et al., 2017; Feldman, 2017), where the sample
 110 complexity lower bound is $\Omega(d^{s^*/2})$, with s^* as the generative exponent (Damian et al., 2024). In
 111 this context, online SGD with batch reusing suffices for learning polynomial link functions (Dandi
 112 et al., 2024; Lee et al., 2024), while for $s^* \geq 3$, only the partial trace estimator proposed by Damian
 113 et al. (2024) can match the SQ lower bound.

114 In the sparse setting, including sparse linear models (Vaskevicius et al., 2019; Zhao et al., 2022;
 115 Gamarnik and Zadik, 2017), sparse PCA (Arous et al., 2020), and planted models (Bandeira et al.,
 116 2022), computational-statistical gaps also exist. Related to our work, Fan et al. (2023) provide a $\tilde{O}(k)$
 117 sample complexity for learning single index models with $q^* = 1$ using diagonal linear networks, and
 118 Neykov et al. (2016) report a $\tilde{O}(k^2)$ result for phase retrieval where $q^* = 2$. However, as previously
 119 noted, the information exponent does not fully characterize the intrinsic computational-statistical

120 tradeoff. Our work completes the picture by providing a gradient-based framework that simultaneously
 121 handles all sparsity levels and any generative exponent $s^* \geq 1$.

122 2 Problem Setup

123 We begin by introducing the notation used in the paper, and then describe the problem setup. For a
 124 probability distribution \mathbb{P} , we denote by $L^2(\mathbb{P})$ the space of square-integrable functions with respect to
 125 \mathbb{P} , and $\stackrel{L^2(\mathbb{P})}{=}$ means equality in $L^2(\mathbb{P})$. We denote the normalized probabilist’s Hermite polynomials by
 126 $\{h_s(\cdot)\}_{s \geq 0}$, where each $h_s(x) := \frac{(-1)^s}{\sqrt{s!}} \cdot e^{x^2/2} \cdot \frac{d^s}{dx^s} e^{-x^2/2}$. These polynomials form an orthonormal
 127 basis for $L^2(\mathcal{N}(0, 1))$, i.e., the space of square-integrable functions under the Gaussian measure.

128 **Gaussian single-index model.** We study the following Gaussian single-index model: The environment
 129 first samples an unobservable signal $\theta^* \sim \pi$ from some known prior $\pi \in \mathcal{P}(\mathbb{S}^{d-1})$. Then i.i.d. data
 130 $(z_1, y_1), \dots, (z_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$ are generated according to the following distribution \mathbb{P}_{θ^*} given θ^* :

$$\mathbb{P}_{\theta^*} : \quad z \sim \mathcal{N}(0, I_d), \quad y \sim p(\cdot | \langle \theta^*, z \rangle). \quad (2.1)$$

131 Here $p(\cdot | \cdot) : \mathbb{R} \mapsto \mathcal{P}(\mathbb{R})$ is referred to as the *link distribution*. A canonical example is the additive
 132 model where $y = p(\langle \theta^*, z \rangle) + \epsilon$ for some deterministic link function $p : \mathbb{R} \rightarrow \mathbb{R}$ (with a slight abuse
 133 of notation) and random noise ϵ . See Damian et al. (2024) for more complicated examples.

134 **Generative exponent.** The following discussion on the generative exponent is based on the work of
 135 Damian et al. (2024). We aim to learn (2.1) where the link distribution p has *generative exponent*
 136 $s^* \geq 1$, a measure of the computational-statistical gap for learning single-index models. We let
 137 $x = \langle \theta^*, z \rangle$. Notice that $\mathbb{P}_{\theta^*}(y, z) = \mathbb{P}(y, x) \cdot \mathcal{N}(z^\perp; 0, I_{d-1})$ where we use \mathbb{P} to denote the joint
 138 distribution of (x, y) as this joint distribution is independent of θ^* . As the marginal distribution
 139 of y is also independent of θ^* , we define the *null distribution* $\mathbb{Q}(y, z) := \mathcal{N}(z; 0, I_d) \otimes \mathbb{Q}(y)$ and
 140 denote $\mathbb{Q}(y, x) := \mathcal{N}(x; 0, 1) \otimes \mathbb{Q}(y)$ where $\mathbb{Q}(y) = \int_{\mathbb{R}} \mathbb{P}(y, x) dx$. It can be shown that under a
 141 square-integrable condition under \mathbb{Q} , the likelihood ratio admits a Hermite expansion with coefficient
 142 functions $\{\zeta_s(y)\}_{s \geq 1}$, i.e.,

$$\frac{\mathbb{P}_{\theta^*}(y, z)}{\mathbb{Q}(y, z)} \stackrel{L^2(\mathbb{Q})}{=} \frac{\mathbb{P}(y, x)}{\mathbb{Q}(y, x)} = \sum_{s=0}^{\infty} \zeta_s(y) \cdot h_s(x), \quad \text{where } \zeta_s(y) = \mathbb{E}_{\mathbb{P}}[h_s(x)|y], \quad (2.2)$$

143 and $\mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \leq 1$ for all $s \geq 1$. Note that (2.2) makes sense only when we are working with the
 144 inner product of \mathbb{P}/\mathbb{Q} and a square-integrable function under the null distribution \mathbb{Q} .

145 **Definition 2.1** (Generative exponent). *For the Gaussian single-index model defined in (2.1), the*
 146 *generative exponent s^* of the link distribution p is defined as $s^*(p) := \min\{s \geq 1 : \mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] > 0\}$.*

147 **Example 2.2** (Example 2.7, Damian et al. (2024)). *Consider the special case of the Gaussian single-*
 148 *index model (2.1) where $y = p(\langle \theta^*, z \rangle)$ for a deterministic link function $p : \mathbb{R} \rightarrow \mathbb{R}$. When p is*
 149 *a polynomial function, it holds that $s^*(p) \leq 2$, and the equality holds if and only if p is an even*
 150 *polynomial. In particular, $s^*(h_s) = 1$ for odd s and $s^*(h_s) = 2$ for even s . While for the example of*
 151 *$p(x) = x^2 \exp(-x^2)$, which is not a polynomial, it has generative exponent $s^*(p) = 4$.*

152 **Two-layer neural networks.** We consider using a two-layer neural network with M hidden neurons
 153 to learn the single-index model (2.1). The weight vector for each neuron $m \in [M]$ is $\theta_m \in \mathbb{R}^d$,
 154 and the weights of the second layer are $a_1, \dots, a_M \in \mathbb{R}$. We collect all the weights and denote
 155 $\boldsymbol{\theta} = (\theta_1, \dots, \theta_M) \in \mathbb{R}^{d \times M}$, $\mathbf{a} = (a_1, \dots, a_M)^\top \in \mathbb{R}^M$. Now for any input $z \in \mathbb{R}^d$, the output of
 156 the network is given by $f(z; \boldsymbol{\theta}, \mathbf{a}) := \sum_{m=1}^M a_m \cdot \sigma(\langle z, \theta_m \rangle)$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is the activation.

157 3 Overview of techniques

158 In this work, we apply gradient-based methods to learn Gaussian single-index models, with a focus
 159 on feature learning in neural networks and the corresponding computational-statistical tradeoff. To
 160 motivate the techniques involved, we begin by discussing an illustrative example that highlights such
 161 tradeoffs. For this overview, we focus on $s^* > 2$ and the uniform prior $\pi = \text{Unif}(\mathbb{S}^{d-1})$. It has been
 162 shown that a gap exists between the information-theoretic lower bound $\Omega(d)$ and the SQ lower bound
 163 $\Omega(d^{s^*/2})$ under this setting when $s^* > 2$ (Bach, 2017; Damian et al., 2024).

164 For illustration, let us consider training a two-layer network with a single neuron under the population
 165 square loss. When the weight of the second layer is small, the reweighted negative gradient g satisfies

$$g = -(2a)^{-1} \nabla_{\theta} (f(z; \theta, a) - y)^2 = -(a \cdot \sigma(\langle z, \theta \rangle) - y) \cdot \sigma'(\langle z, \theta \rangle) \cdot z \approx y \sigma'(\langle z, \theta \rangle) \cdot z.$$

166 Taking expectation over $(z, y) \sim \mathbb{P}_{\theta^*}$ and using the likelihood ratio decomposition in (2.2), we have

$$\mathbb{E}_{\mathbb{P}_{\theta^*}} [g] \approx \underbrace{\mathbb{E}_{\mathbb{Q}} [y] \cdot \mathbb{E}_{\mathbb{Q}} [\sigma'(\langle z, \theta \rangle) \cdot z]}_{\text{bias}} + \sum_{s \geq s^*} \underbrace{\mathbb{E}_{\mathbb{Q}} [y \zeta_s(y)] \cdot \mathbb{E}_{\mathbb{Q}} [h_s(\langle \theta^*, z \rangle) \cdot \sigma'(\langle z, \theta \rangle) \cdot z]}_{\text{informative queries}}, \quad (3.1)$$

167 where we use the fact that y and z are independent under the null distribution \mathbb{Q} . Note that the
 168 *bias* term does not contain any information about θ^* , and it can be easily removed by a debiasing
 169 procedure, so we assume for simplicity that $\mathbb{E}[y] = 0$.

170 **Failure of vanilla online minibatch SGD** We first consider the vanilla online minibatch SGD, which
 171 updates the weight vector θ by $\theta \leftarrow \theta - \eta \sum_{i=1}^n g_i$ for a minibatch of size n . The sample complexity of
 172 gradient-based methods is determined by the signal-to-noise ratio (SNR) of the one-sample gradient,
 173 which in our case is defined as $\text{SNR} := \mathbb{E}[\langle g, \theta^* \rangle]^2 / \mathbb{E}[\|g\|_2^2]$. This is the square of the alignment
 174 between g and θ^* , governed primarily by the informative query corresponding to the lowest degree
 175 s^* in (3.1) assuming that $\mathbb{E}_{\mathbb{Q}} [y \zeta_{s^*}(y)] \neq 0$. It can be shown that the inner product between the
 176 lowest-degree informative query in (3.1) and the signal θ^* satisfies (see Lemma H.1)

$$\mathbb{E}_{\mathbb{Q}} [h_{s^*}(\langle \theta^*, z \rangle) \cdot \sigma'(\langle z, \theta \rangle) \cdot \langle z, \theta^* \rangle] \approx s^* \cdot \hat{\sigma}_{s^*} \cdot \langle \theta^*, \theta \rangle^{s^*-1} = \hat{\sigma}_{s^*} \cdot O(d^{-(s^*-1)/2}), \quad (3.2)$$

177 where $\hat{\sigma}_{s^*}$ is the s^* -th coefficient in the Hermite expansion of σ . While for $\|g\|_2$, we have

$$\mathbb{E}_{\mathbb{P}_{\theta^*}} [\|g\|_2^2] \approx d \cdot \mathbb{E}_{\mathbb{Q}} [y^2 \sigma'(\langle z, \theta \rangle)^2] = \Omega(d),$$

178 where the high-order terms in the likelihood ratio decomposition are ignored and we come back to
 179 this point later. Now we can argue why vanilla online minibatch SGD has difficulty achieving the
 180 SQ lower bound for generative exponent $s^* > 2$: Suppose $\mathbb{E}_{\mathbb{Q}} [y \zeta_{s^*}(y)]$ and $\hat{\sigma}_{s^*}$ are both nonzero
 181 constants. Then the one-sample SNR is $O(d^{-s^*})$. For a minibatch with n samples, the SNR of the
 182 gradient averaged over the minibatch is roughly n times the one-sample SNR², i.e., nd^{-s^*} . To ensure
 183 one update step achieves alignment, i.e., the square root of the n -sample SNR, $\sqrt{nd^{-s^*}}$, exceeding
 184 the trivial $d^{-1/2}$ threshold attained by a random vector, it requires at least d^{s^*-1} samples. Note that
 185 the sample complexity would become even worse if $s^* < \arg \min_{s \geq s^*} \{s : \mathbb{E}_{\mathbb{Q}} [y \zeta_s(y)] \neq 0\}$. This
 186 contrasts with the sample complexity $O(d^{s^*/2})$ suggested by the SQ lower bound.

187 The above failure of vanilla online minibatch SGD exposes three key challenges:

- 188 (i) **(Non-polynomial)** How to handle the infinite sum of high-order terms in the likelihood ratio?
- 189 (ii) **(Low SNR)** How to enhance the SNR to achieve the SQ lower bound?
- 190 (iii) **(Zero correlation)** How to ensure that the algorithm still works if $\mathbb{E}_{\mathbb{Q}} [y \zeta_{s^*}(y)] = 0$?

191 Below we discuss our techniques for addressing these challenges.

192 **Label transformation via general gradient oracle.** The idea to fix the zero correlation problem
 193 is to apply a nonlinear transformation $\mathcal{T} : \mathbb{R} \rightarrow \mathbb{R}$ to y such that $\mathcal{T}(y)$ has nonzero correlation with
 194 $\zeta_{s^*}(y)$. This label transformation technique has been widely used in the literature (Lu and Li, 2020;
 195 Mondelli and Montanari, 2018; Dudeja and Hsu, 2018; Chen and Meka, 2020; Damian et al., 2024).
 196 In particular, Lee et al. (2024) show that the label transformation can be automatically realized by
 197 running two gradient steps on the same batch, a technique termed as *batch-reusing* (Dandi et al.,
 198 2024; Arnaboldi et al., 2024). In this work, we study a more *general class of gradient-based methods*
 199 with gradient of form $g = \psi(y, \langle \theta, z \rangle) z$, which is an abstract form of the transformed gradient
 200 $\mathcal{T}(y) \sigma'(\langle z, \theta \rangle) z$. The desired condition becomes $\mathbb{E}_{\mathbb{Q}} [\hat{\psi}_{s^*-1}(y) \zeta_{s^*}(y)] \neq 0$, where $\hat{\psi}_s(y)$ is the s -th
 201 Fourier coefficient function of $\psi(y, x)$ in the Hermite basis of x . One particular way to obtain such a
 202 gradient is to use a modified loss function, similar to the approach in Joshi et al. (2024), while in our
 203 case the specific choice of ψ is also related to the other two challenges addressed as follows.

204 **Exploration by weight perturbation with high-pass activation.** The low-SNR challenge corre-
 205 sponds to the fact that points on the equator of \mathbb{S}^{d-1} orthogonal to θ^* are all saddle points in terms of

²This argument is not fully rigorous because $\mathbb{E}_{\mathbb{P}_{\theta^*}} [\|g\|_2^2]$ also includes ‘‘bias’’ $\|\mathbb{E}_{\mathbb{P}_{\theta^*}} [g]\|_2^2$ besides the fluctuations, but it remains valid as long as $\|g\|_2^2$ is dominated by fluctuations from all d directions at initialization.

Algorithm 1 Gradient-based Feature Learning for Uniform Signal Prior

- 1: **Input:** Initialization $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_M^{(0)}) \in \mathbb{R}^{d \times M}$, where $\theta_m^{(0)} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1})$, $\mathbf{a} = a \cdot \mathbf{1} \in \mathbb{R}^M$, number of iterations $T \in \mathbb{N}$, learning rate $\eta > 0$, batch size $n \in \mathbb{N}$, polarization level $\gamma \in (0, 1)$, number of perturbation $L \in \mathbb{N}$.
 - 2: **for** iteration $t = 0, 1, \dots, T - 1$ **do**
 - 3: Sample a fresh mini-batch of data $\{(z_i^{(t)}, y_i^{(t)})\}_{i=1}^n$.
 - 4: Perturb weights $w_{m,l}^{(t)} = (\gamma\theta_m^{(t)} + \xi_{m,l}^{(t)}) / \|\gamma\theta_m^{(t)} + \xi_{m,l}^{(t)}\|_2$, $\xi_{m,l}^{(t)} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1})$ for all m, l .
 - 5: Compute the gradients $g_{m,l,i}^{(t)} = (\psi(y_i^{(t)}, \langle w_{m,l}^{(t)}, z_i^{(t)} \rangle) + \text{err}_{m,l,i}^{(t)}) \cdot z_i^{(t)}$ for all m, l, i .
 - 6: Aggregate the gradients: $g_m^{(t)} = (nL)^{-1} \sum_{i=1}^n \sum_{l=1}^L (g_{m,l,i}^{(t)} - \hat{\psi}_1(y_i^{(t)})w_{m,l}^{(t)})$ for all m .
 - 7: Normalize the update step: $\bar{g}_m^{(t)} = g_m^{(t)} / \|g_m^{(t)}\|_2$ for all m .
 - 8: Update the weights $\theta_m^{(t+1)} = (\theta_m^{(t)} + \eta\bar{g}_m^{(t)}) / \|\theta_m^{(t)} + \eta\bar{g}_m^{(t)}\|_2$ for all m .
 - 9: **end for**
 - 10: **Output:** Final model weights $\theta^{(T)}$.
-

206 $|\langle \theta, \theta^* \rangle|$, and random initialization typically lies near this equator. To efficiently escape from such
 207 saddle points, we perform random weight perturbation, akin to the approach in Jin et al. (2017) for
 208 non-convex optimization. Specifically, suppose the activation σ is high-pass and has the lowest degree
 209 s^* , i.e., $\sigma(x) = \sum_{s \geq s^*} \hat{\sigma}_s h_s(x)$, and consider for simplicity the case of odd s^* . In the extreme case
 210 where θ is perturbed into i.i.d. pure noise $\theta_1, \dots, \theta_L \sim \text{Unif}(\mathbb{S}^{d-1})$, we compute the gradient for
 211 each θ_l and aggregate them into $g = L^{-1}(g_1 + \dots + g_L)$. Using the properties of the Gaussian noise
 212 operator (see Appendix B for details), the second moment of this aggregated gradient satisfies

$$\mathbb{E}[\|g\|_2^2] \approx \frac{d}{L^2} \sum_{l,l'=1}^L \mathbb{E}_{\mathbb{Q}}[y^2] \cdot \mathbb{E}_{\mathbb{Q}}[\sigma'(\langle z, \theta_l \rangle) \sigma'(\langle z, \theta_{l'} \rangle)] \approx d \sum_{s \geq s^*} s \cdot \hat{\sigma}_s^2 \cdot \mathbb{E}_{\theta, \theta'}[\langle \theta, \theta' \rangle^{s-1}],$$

213 where θ, θ' are drawn independently from $\text{Unif}(\mathbb{S}^{d-1})$. Since $\langle \theta, \theta' \rangle \approx d^{-1/2}$, we have $\mathbb{E}[\|g\|_2^2] \approx$
 214 $O(d^{-(s^*-3)/2})$, yielding a higher one-sample SNR as the first moment remains unchanged and
 215 pushing the sample complexity towards the SQ lower bound. Moreover, we also see from the
 216 above calculation that the weight perturbation resolves the non-polynomial issue thanks to the near-
 217 orthogonality of the perturbed weights. The above heuristics can be made rigorous for polynomially
 218 large L , thereby handling non-polynomial link and activation functions.

219 Our approach also draws inspiration from the landscape smoothing method in Damian et al. (2024),
 220 but in contrast to their problem setup, we do not require full knowledge of the link distribution in
 221 advance. Instead, it suffices to know the generative exponent s^* to construct a high-pass activation
 222 function as well as the gradient oracle ψ . See Example 4.6 for a detailed discussion on this.

223 4 Gradient-based Algorithm for Uniform Prior

224 We first present our method and results for the case of $\theta^* \sim \text{Unif}(\mathbb{S}^{d-1})$, or equivalently, when
 225 there is no structural information on θ^* . Motivated by the discussion in Section 3, we propose a
 226 gradient-based algorithm (Algorithm 1) that can train a two-layer neural network to learn the unknown
 227 signal θ^* with $\tilde{O}(d^{s^*/2} \vee d)$ sample complexity, nearly matching the corresponding SQ lower bound.

228 4.1 Gradient-based Training Algorithm (Algorithm 1)

229 We initialize each neuron m with $\theta_m^{(0)} \sim \text{Unif}(\mathbb{S}^{d-1})$, and we set $a_m^{(t)} \equiv a$ for some sufficiently small
 230 $a > 0$ throughout the training. In each iteration $t \in [T]$, we sample a new data batch of size n .

231 **Weight perturbation.** Before calculating the gradients, we first perturb the weights of each neuron to
 232 get L noisy replica, by injecting uniform noise from the sphere \mathbb{S}^{d-1} as in Line 4. There is a simple
 233 rule for choosing the polarization level γ . In the previous section, we discussed how $\mathbb{E}_{\mathbb{P}_{\theta^*}}[\|g\|_2^2]$ in
 234 the one-sample SNR depends on the following quantity:

$$\mathbb{E}_{l,l'} \langle w_{m,l}^{(t)}, w_{m,l'}^{(t)} \rangle^{s^*-1} \lesssim (\gamma^2 \|\theta_m^{(t)}\|_2^2)^{s^*-1} + \mathbb{E}_{l,l'} \langle \xi_{m,l}^{(t)}, \xi_{m,l'}^{(t)} \rangle^{s^*-1} \approx (\gamma^2 \vee d^{-1/2})^{s^*-1}.$$

235 In this context, γ^2 represents the bias from the *exploitation* of the learned search direction, and $d^{-1/2}$
 236 accounts for the variance from the *exploration* for the unknown signal. In fact, γ should be set as
 237 large as possible to maximize exploitation while still ensuring that the exploration noise dominates.
 238 This balance is necessary to fully gain the SNR enhancement from weight perturbation. This gives
 239 rise to the choice $\gamma = \tilde{\Theta}(d^{-1/4})$. Moreover, it suffices to set $L = \tilde{\Omega}(n/\sqrt{d})$ as stated in Theorem 4.2.

240 **Gradient aggregation and debiasing.** Then for each neuron m and its perturbed weights $w_{m,l}^{(t)}$, we
 241 calculate the gradient $g_{m,l,i}^{(t)}$ for every sample (Line 5). The gradient is expressed by decoupling the
 242 primary search direction $\psi(y_i, \langle w_{m,l}, z_i \rangle) \cdot z_i$ from the error term $\text{err}_{m,l,i} \cdot z_i$ (omitting the time
 243 index t). For two-layer neural networks, we discussed in the previous section an example where
 244 $\psi(y_i, \langle w_{m,l}, z_i \rangle) = \mathcal{T}(y_i)\sigma'(\langle w_{m,l}, z_i \rangle)$ for some transformation \mathcal{T} . While the error term, arising
 245 from neuron interactions during backpropagation, can be reduced by keeping the weights of the
 246 second layer sufficiently small. We provide the assumptions and more examples of ψ in Section 4.2.

247 Next, we aggregate the gradients for each neuron m by averaging over the n samples and L perturba-
 248 tions to get $g_m^{(t)}$ as shown in Line 6. Here we additionally subtract a term $\hat{\psi}_1(y_i^{(t)})w_{m,l}^{(t)}$ to debias the
 249 gradient, where $\hat{\psi}_1(y)$ is the first coefficient function in the Hermite expansion of the oracle function
 250 $\psi(y, x)$ with respect to x . Finally, we update $\theta_m^{(t)}$ according to Line 7 and Line 8.

251 4.2 Feature Alignment and Statistical Complexity

252 **Assumption 4.1.** For the Gaussian single-index model in (2.1) with generative exponent $s^* \geq 1$, the
 253 oracle $\psi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ satisfies the following conditions:

- 254 (a) (Quadruple-integrable under \mathbb{Q}). Both $\psi(y, x)^2 x^2$ and $\psi(y, x)^2$ are square-integrable under
 255 the null distribution \mathbb{Q} . Therefore, $\psi \in L^2(\mathbb{Q})$ admits an Hermite expansion $\psi(y, x) \stackrel{L^2(\mathbb{Q})}{=} \sum_{s=0}^{\infty} \hat{\psi}_s(y) \cdot h_s(x)$, where $\hat{\psi}_s(y)$ is the s -th coefficient function and $\sum_{s=0}^{\infty} \mathbb{E}_{\mathbb{Q}}[\hat{\psi}_s(y)^2] < \infty$.
 256 (b) (High-pass under \mathbb{Q}). For all $s = 1, \dots, s^* - 2$, the s -th coefficient function is zero, i.e., $\hat{\psi}_s(y) \equiv 0$.
 257 In addition, there exists a constant $C > 0$ such that $|\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)]| \geq C$.
 258 (c) (Polynomial-like tail under \mathbb{P} and \mathbb{Q}). There exists a constant $C > 0$ such that for all $r \geq 1$,
 259 $\max\{\mathbb{E}_{\mathbb{P}}[|\psi(y, x)|^r], \mathbb{E}_{\mathbb{Q}}[|\psi(y, x)|^r]\} \leq C \cdot r^{Cr}$.
 260

261 The quadruple-integrability condition (Assumption 4.1(a)) ensures that the decomposition of the
 262 likelihood ratio in (2.2) is well defined for calculations involving the second moment, i.e.,

$$\mathbb{E}_{\mathbb{P}}[\psi(y, x)^2 x^2] = \mathbb{E}_{\mathbb{Q}}\left[\psi(y, x)^2 x^2 \cdot \frac{\mathbb{P}(y, x)}{\mathbb{Q}(y, x)}\right] = \mathbb{E}_{\mathbb{Q}}\left[\psi(y, x)^2 x^2 \cdot \sum_{s=0}^{\infty} \zeta_s(y) h_s(x)\right].$$

263 The high-pass condition (Assumption 4.1(b)) has been motivated in Section 3 and guarantees noise
 264 reduction for the second moment. The polynomial-like tail condition (Assumption 4.1(c)) is used
 265 for concentration arguments in the proof. Note that this condition is analogue to the Gaussian
 266 hypercontractivity property, where $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[|f(x)|^r] \lesssim r^{Dr/2}$ if $f(x)$ is a polynomial of degree at
 267 most D . In particular, ψ can be constructed as $\psi(y, x) = \ell'(y)\sigma'(x)$, where ℓ is the loss function and
 268 σ is the activation in the two-layer network. It suffices to use a loss ℓ with bounded derivative and a
 269 polynomial activation σ for the polynomial-like tail condition (see Section 4.2.1 and 4.2.2).

270 Now we are ready to state our first main result on the sample complexity of Algorithm 1 for uniform
 271 prior. See Appendix D for a proof sketch of the theorem and Appendix E for a detailed proof.

272 **Theorem 4.2** (Sample complexity for uniform prior). Under Assumption 4.1, set the initializa-
 273 tion of the weights as $\theta_1^{(0)}, \dots, \theta_M^{(0)} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1})$. Suppose that the event $\mathcal{E} = \{|\text{err}_{m,l,i}^{(t)}| \leq$
 274 $d^{-10s^*}, \forall (m, l, i, t) \in [M] \times [L] \times [n] \times [T]\}$ holds with probability at least $1 - O(d^{-c_0})$ for some
 275 constant $c_0 > 0$ with (M, L, n, T) specified as follows. Set the learning rate $\eta > 2$, the polarization
 276 level $\gamma = d^{-1/4}(\log d)^{1/4}$, and the number of neurons $M = \Theta(1)$. Suppose

$$n = \Theta((d^{s^*/2}(\log d)^{1+s^*/2}) \vee d(\log d)^2), \quad L = \Theta(d^{(s^*+1)/2} \log d).$$

277 Define $\tau = \eta^{-2}/(1 - \eta^{-1})^2$, and let $\Delta = (\log d)^{-1/2}$ if $s^* \leq 2$ and $\Delta = d^{-1/4}(\log d)^{1/4}$ if $s^* \geq 3$.
 278 Then with probability at least $1 - O(d^{-c})$ for some constant $c > 0$, after running Algorithm 1 for

279 $T = O(\log d + \log(\Delta^{-1})/\log(\tau^{-1}))$ steps, there are at least $\Omega(M)$ neurons having alignment
 280 $|\langle \theta_m^{(T)}, \theta^* \rangle| \geq 1 - O(\sqrt{\Delta})$.

281 Theorem 4.2 shows that the sample complexity of Algorithm 1 is $nT = \tilde{\Theta}(d^{s^*/2} \vee d)$, matching the
 282 SQ sample complexity lower bound for all $s^* \geq 1$ established by Damian et al. (2024). Compared to
 283 the partial trace method in Damian et al. (2024), our algorithm does not require special warm-start
 284 initialization. Meanwhile, the computational complexity of Algorithm 1 is $MLnT = \tilde{\Theta}(d^{s^*+1/2} \vee d^2)$.
 285 So far the gradient oracle ψ is still an abstract object, and next we will instantiate the above general
 286 theorem with concrete examples of ψ that yield implementable algorithms.

287 **Remark 4.3** (Benefit of overparametrization). *Algorithm 1 trains a two-layer neural network with*
 288 *constant width M , involving L times of perturbation for every neuron in each step. Indeed, this is*
 289 *equivalent to train a two-layer neural network with width $LM = \Theta(d^{(s^*-1)/2} \vee (d \log d)^{1/2})$, where*
 290 *we divide the neurons into L groups, each having M neurons. In each iteration we perturb the weights*
 291 *and compute the gradients, and then aggregate the gradients within each group of M neurons. This*
 292 *combination of weight perturbation and gradient sharing exploits the benefit of overparametrization.*

293 4.2.1 Online SGD with Batch Reusing

294 The oracle function ψ can be specialized to two consecutive gradient descent steps on the same batch.

295 **Example 4.4** (Batch-reusing: ψ for polynomial link functions). *Suppose that the link distribution is*
 296 *a polynomial of degree q . We consider ψ induced by batch-reusing on single neuron, i.e., $\psi(y, x) =$
 297 $y\sigma'(x) + y\sigma'(x) + y\sigma'(x)$ (see Section 4.2 of Lee et al. (2024) for deduction) and choose $\sigma'(x) =$
 298 $\sum_{i=0}^q c_i h_i(x)$ where $C_q \in \mathbb{N}$ is a constant depending only on q and each $c_i \sim \text{Unif}([0, 1])$.*

299 **Corollary 4.5** (Batch-reusing for polynomial link function). *Suppose that the link distribution is*
 300 *given by a polynomial link function. Under the same setups in Theorem 4.2 with the oracle ψ given by*
 301 *Example 4.4, the sample complexity of Algorithm 1 is $\tilde{\Theta}(d)$, recovering the result of Lee et al. (2024).*

302 The proof is deferred to Appendix C.2.1. However, batch-reusing may not be optimal for $s^* \geq 3$ due
 303 to violation of the high-pass condition, necessitating a more general approach to construct ψ .

304 4.2.2 Label Transformation via Modified Loss

305 We discuss another approach to construct ψ by modifying the loss function, a universal method for
 306 arbitrary generative exponent $s^* \geq 1$. Additional details and proofs are postponed to Appendix C.2.2.

307 **Example 4.6** (ψ based on modified loss). *Let $\psi(y, x) = \ell'(y)\sigma'(x)$ with $\ell(y)$ being certain loss*
 308 *function and $\sigma(x)$ being some activation function. Such a form corresponds to the gradient of the*
 309 *loss $\ell(y - f(z; \theta))$ (assuming that \mathbf{a} is fixed and has small entries), since*

$$a_m^{-1} \nabla_{\theta_m} \ell(y - f(z; \theta)) = -\ell'(y - f(z; \theta)) \cdot \sigma'(\langle \theta_m, z \rangle) \cdot z = \underbrace{-\ell'(y) \cdot \sigma'(\langle \theta_m, z \rangle)}_{:= \psi(y, \langle \theta_m, z \rangle)} \cdot z + \text{err}_m \cdot z,$$

310 where $\text{err}_m := [\ell'(y) - \ell'(y - f(z; \theta))] \cdot \sigma'(\langle \theta_m, z \rangle) = O(f(z; \theta, \mathbf{a})) \cdot \sigma'(\langle \theta_m, z \rangle)$ denotes small
 311 error for sufficiently small \mathbf{a} . In Appendix C.2.2, we provide specific choice of the activation function
 312 $\sigma(x)$ (order- s^* Hermite polynomial) and the loss function $\ell(y)$ (a carefully designed random loss
 313 function), satisfying all the conditions in Assumption 4.1.

314 **Corollary 4.7** (Modified loss for general s^*). *The oracle ψ given by Example 4.6 satisfies all the*
 315 *assumptions of Theorem 4.2, thus the results of Theorem 4.2 hold for Algorithm 1 using this ψ .*

316 5 Exploiting the Structure: Algorithm for Sparse Prior

317 We have seen that the sample complexity scales polynomially with the ambient dimension d when
 318 the prior on θ^* is uninformative, and our method achieves nearly optimal computational-statistical
 319 tradeoff under the SQ framework. It is natural to ask whether our method can benefit from extra
 320 structural information on θ^* , one classic example being sparsity. Towards this end, we consider an
 321 extension of the framework in the previous section to the setting where θ^* is a k -sparse vector.

322 **Gaussian single-index model with sparse signal.** Given sparsity level $k = o(\sqrt{d})$, we consider the
 323 Gaussian single-index model in (2.1) with θ^* drawn from a k -sparse prior:

$$\pi_k : \quad \theta^* | \phi^* \sim \text{Unif}(\mathbb{S}^{k-1}(\phi^*)), \quad \phi^* \sim \text{Unif}(\mathcal{S}_k), \quad (5.1)$$

324 where $\mathcal{S}_k := \{\phi \subset [d] : |\phi| = k\}$ is the collection of all k -sparse support sets, and $\mathbb{S}^{k-1}(\phi) := \{x \in$
 325 $\mathbb{R}^d : \sum_{i \in \phi} x_i^2 = 1, x_j = 0, \forall j \notin \phi\}$ is the associated k -dimensional unit sphere for any $\phi \in \mathcal{S}_k$.

326 5.1 Algorithm design: How to leverage sparsity?

327 Note that Algorithm 1 can also learn the k -sparse Gaussian single-index model, albeit with $\tilde{O}(d^{s^*/2} \vee$
 328 $d)$ samples, which is apparently suboptimal in light of the classic example of sparse linear regression.
 329 Here the key challenge is *support identification* of ϕ^* , and the issue of Algorithm 1 lies in the weight
 330 perturbation using noise $\xi \sim \text{Unif}(\mathbb{S}^{d-1})$, thus unaware of the sparsity of θ^* . Below we discuss how
 331 to calibrate the weight perturbation with the sparse prior.

332 **Perturbation by replicating the prior is not enough.** An intuitive first-cut attempt is to use perturba-
 333 tion noise from the same distribution as the sparse prior π_k in (5.1), which turns out to be suboptimal
 334 as well. To illustrate this, we assume for simplicity a balanced θ^* where every nonzero entry of θ^* is
 335 equal to $k^{-1/2}$, and consider i.i.d. $\xi_1, \dots, \xi_L \sim \pi_k$. For each $j \in \phi^*$, consider the j -th entry of the
 336 lowest-degree informative query (analogous to (3.2)), whose first moment satisfies

$$\mathbb{E}_{\mathbb{P}_{\theta^*}}[g_j] \approx \mathbb{E}_{\mathbb{Q}}[y \zeta_s(y)] \cdot s^* \hat{\sigma}_{s^*} \cdot \frac{1}{L} \sum_{l=1}^L \langle \theta^*, \theta_l \rangle^{s^*-1} \theta_j \approx \frac{\mathbb{E}_{\theta \sim \pi_k}[\langle \theta^*, \theta \rangle^{s^*-1}]}{\sqrt{k}} \simeq \frac{k^2}{d} \cdot \frac{k^{-(s^*-1)}}{\sqrt{k}},$$

337 where the last step follows from direct calculation for $\theta \sim \pi_k$. Similarly, the second moment satisfies

$$\mathbb{E}_{\mathbb{P}_{\theta^*}}[g_j^2] \approx \sum_{s \geq s^*} s \cdot (\hat{\sigma}_s)^2 \cdot \mathbb{E}_{\theta, \theta' \sim \pi_k}[\langle \theta, \theta' \rangle^{s-1}] \simeq \frac{k^2}{d} \cdot k^{-(s-1)},$$

338 where θ and θ' are drawn independently from the prior π_k . This calculation implies that the fluctuation
 339 of each entry of the aggregated gradient is of order $\sqrt{k^2/d} \cdot k^{-(s^*-1)/2} \cdot n^{-1/2}$ for batch size n . To
 340 successfully identify the true support ϕ^* , the signal must be larger than the fluctuation for entries
 341 in ϕ^* , resulting in a sample complexity of $n = \tilde{O}(k^{s^*} \cdot d/k^2)$. In comparison to the SQ lower
 342 bound in Theorem 5.4, this is suboptimal by a factor d/k^2 . However, for $s^* = 1$, we observe that
 343 $\mathbb{E}_{\mathbb{P}_{\theta^*}}[g_j] = \Omega(k^{-1})$ for $j \in \phi^*$ and $\mathbb{E}_{\mathbb{P}_{\theta^*}}[g_j^2] = O(1)$, indicating that the support ϕ^* can still be
 344 identified using $\tilde{O}(k)$ samples. The form of the perturbation does not matter here since both $\langle \theta^*, \theta \rangle$
 345 and $\langle \theta, \theta' \rangle$ are degree zero in terms of k^{-1} . Therefore, we conjecture that our algorithm succeeds for
 346 $s^* = 1$ even without weight perturbation as outlined in Conjecture 5.3.

347 **Perturbation by groups that cover the prior.** The suboptimality of the previous strategy originates
 348 from the fact that ϕ^* is sampled from a uniform distribution over $\binom{d}{k}$ different k -sparse support sets,
 349 making it unlikely for two independent sets to overlap (only with k^2/d probability). Then how to
 350 perturb the weights in a way that guarantees a significant overlap with ϕ^* ? The solution is to construct
 351 a *polynomial-size cover* for the prior π_k . Specifically, we divide \mathcal{S}_k into d subsets, where the j -th
 352 subset is define as $\mathcal{S}_{k,j} := \{\phi \in \mathcal{S}_k \mid j \in \phi\}$, which contains all k -sparse support sets that include
 353 the j -th coordinate. Now suppose $\theta^* \in \mathcal{S}_{k,j}$, then for any perturbed weight θ_l with support from the
 354 same subset $\mathcal{S}_{k,j}$, its support overlaps with ϕ^* almost surely, thereby eliminating the d/k^2 factor.

355 In particular, considering a two-layer neural network with width d , the above strategy can be carried
 356 out by perturbing each neuron m using $\theta_{m,1}, \dots, \theta_{m,L}$ whose support sets are sampled from the
 357 same group $\mathcal{S}_{k,m}$. As a result, at least k neurons will have one or more overlapping coordinates with
 358 the true signal θ^* . For these neurons, the signal in the aggregated gradient would be strong enough
 359 for simple thresholding methods to correctly identify the true support ϕ^* with $\tilde{O}(k^*)$ samples. We
 360 further refine this by first projecting the aggregated gradient for each neuron onto its top- k support,
 361 and then selecting the *strongest* projected gradient to update the weights.

362 Combining these yields Algorithm 2 for the sparse case, where we define the support projection
 363 matrix $P_\phi := \sum_{i \in \phi} e_i e_i^\top$ and the top- k operator $\text{Top}_k(v) := \text{argmax}_{\phi \subset \mathcal{S}_k} \|P_\phi(v)\|_1$, which extracts
 364 the k -sparse support ϕ corresponding to the largest (in absolute value) k entries of v . We set the
 365 polarization level $\gamma = k^{-1/2}$, following the same balance between exploitation and exploration as in
 366 the uniform case, since the exploration noise is now of order k^{-1} .

367 5.2 Sample Complexity Analysis for Sparse Prior

368 **Theorem 5.1** (Sample complexity for sparse prior). *Under Assumption 4.1, consider the sparse prior*
 369 *in (5.1) with sparsity level k satisfying $\omega(d^\iota) < k < o(\sqrt{d})$ for a small $\iota > 0$. Suppose that the*

Algorithm 2 Gradient-based Feature Learning for Sparse Signal Prior

- 1: **Input:** Initialization $\theta^{(0)} = (\theta_1^{(0)}, \dots, \theta_M^{(0)}) \in \mathbb{R}^{d \times M}$, where $\theta_m^{(0)} = e_m$, $\mathbf{a} = a \cdot \mathbf{1} \in \mathbb{R}^M$ with number of neurons $M = d$, number of iterations $T \in \mathbb{N}$, batch size $n \in \mathbb{N}$, polarization level $\gamma \in (0, 1)$, number of perturbation $L \in \mathbb{N}$.
 - 2: **for** iteration $t = 0, 1, \dots, T - 1$ **do**
 - 3: Sample a fresh mini-batch $\{(z_i^{(t)}, y_i^{(t)})\}_{i=1}^n$.
 - 4: Perturb as Line 4 in Algorithm 1 with $\xi_{m,l}^{(t)} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{k-1}(\phi_{m,l}))$ and $\phi_{m,l} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathcal{S}_{k,m})$.
 - 5: Compute and aggregate the gradients to get $g_m^{(t)}$, same as Line 5&6 in Algorithm 1.
 - 6: Find the top- k support of $g_m^{(t)}$ and project: $\phi_m^{(t)} = \text{Top}_k(g_m^{(t)})$, $\tilde{g}_m^{(t)} = P_{\phi_m^{(t)}}(g_m^{(t)})$ for all m .
 - 7: Locate the neuron with the largest $\|\tilde{g}_m^{(t)}\|_2$: $\hat{m} = \text{argmax}_m \|\tilde{g}_m^{(t)}\|_2$.
 - 8: Update weights by gradient sharing: $\theta_m^{(t+1)} = \tilde{g}_m^{(t)} / \|\tilde{g}_m^{(t)}\|_2$ for all m .
 - 9: **end for**
 - 10: **Output:** Model weights $\theta^{(T)}$.
-

370 event $\mathcal{E} = \{|\text{err}_{m,l,i}^{(t)}| \leq d^{-10s^*}, \forall (m, l, i, t) \in [M] \times [L] \times [n] \times [T]\}$ holds with probability at least
 371 $1 - O(d^{-c_0})$ for some constant $c_0 > 0$ with (M, L, n, T) specified as follows. Let $\gamma = k^{-1/2}$ and

$$n = \Omega((k \log^3 k)^{s^*} \cdot \log d), \quad L = \Omega(k^{(s^*+3)/2} \cdot \log(k)^{s^*-1}).$$

372 Then with probability at least $1 - O(k^{-c})$ for some $c > 0$, after running Algorithm 2 with $T = 2$
 373 iterations, there are at least $\Omega(M)$ neurons that have alignment $|\langle \theta_m^{(T)}, \theta^* \rangle| \geq 1 - O(\Delta)$, where
 374 $\Delta = k^{-1} \vee (k^{-(s^*-1)/2} \cdot \log(k)^{-3/2}) = o(1)$.

375 Theorem 5.1 shows that the sample complexity of Algorithm 2 is $n = \tilde{O}(k^{s^*})$, matching the SQ lower
 376 bound established in Theorem 5.4 which will be presented below. Here for simplicity, we essentially
 377 use an infinitely large learning rate when updating the weights in Line 8 of Algorithm 2, so it takes
 378 only two iterations to achieve strong alignment. This is equivalent to running the same algorithm with
 379 a finite learning rate but with a larger number of iterations, which is omitted for brevity.

380 **Remark 5.2** (Implication for sparse tensor PCA). *The connection between the Gaussian single-index*
 381 *model and tensor PCA has been discussed in Damian et al. (2023), by showing that estimating*
 382 *θ^* corresponds to a tensor PCA problem defined over the empirical Hermite tensors. Our weight*
 383 *perturbation technique can be potentially applied to iteratively solve sparse tensor PCA problems.*

384 Next we present the conjecture on the success of our algorithm for $s^* = 1$ discussed in Section 5.1.

385 **Conjecture 5.3.** *For $d^t < k < o(d)$ with $s^* = 1$, Algorithm 2 succeeds with sample complexity*
 386 *$n = \tilde{O}(k)$. Furthermore, the same guarantee applies even without perturbing the weights.*

387 Finally, we present the following SQ lower bound for the sparse prior, complementing Theorem 5.1.

388 **Theorem 5.4** (SQ lower bound). *Consider the Gaussian single-index model in (2.1) with generative*
 389 *exponent $s^* \geq 1$. Suppose θ^* is k -sparse for $\omega((\log d)^2) \leq k \leq d/2$. Take $c > 2$ as a constant. For*
 390 *any (stochastic) algorithm using $\exp(\Omega((\log d)^c))$ calls to the $\text{VSTAT}(\mathbb{P}_{\theta^*}, n)$ oracle with sample*
 391 *size n , in order to achieve nontrivial alignment $|\langle \hat{\theta}, \theta^* \rangle| > \rho$ with probability at least $2/3$, it requires*

$$n \gtrsim \frac{k^{s^*}}{(\log d)^{cs^*}}, \quad \text{where } \rho = \tilde{\omega}(k^{-1}) \quad \text{if } (\log d)^2 < k < \sqrt{d(\log d)^c}, \quad (5.2)$$

$$n \gtrsim \frac{d^{s^*/2}}{(\log d)^{cs^*/2}}, \quad \text{where } \rho = \tilde{\omega}(d^{-1/2}) \quad \text{if } \sqrt{d(\log d)^c} \leq k \leq d/2. \quad (5.3)$$

392 In fact, the effective SQ lower bound should be no smaller than the information-theoretic lower
 393 bound $\Omega(k \log(d/k))$ (Neykov et al., 2016). For $k = o(\sqrt{d})$, running Algorithm 2 will succeed with
 394 $\tilde{O}(k^{s^*})$ samples, matching the lower bound in (5.2) for every $s^* \geq 1$. For $k = \Omega(\sqrt{d})$, running
 395 Algorithm 1 will succeed with $\tilde{O}(d^{s^*/2})$ samples, matching the lower bound in (5.3) for $s^* \geq 2$. In
 396 addition, for $k = \Omega(\sqrt{d})$ with $s^* = 1$, we conjecture $n = \tilde{O}(k)$ samples to be sufficient, where the
 397 information-theoretic lower bound is $\Omega(k \log(d/k))$. This gives rise to the paradigm in Figure 1 (b).

398 **References**

- 399 Abbe, E., Adsera, E. B. and Misiakiewicz, T. (2023). Sgd learning on neural networks: leap com-
400 plexity and saddle-to-saddle dynamics. In *The Thirty Sixth Annual Conference on Learning Theory*.
401 PMLR.
- 402 Allen-Zhu, Z. and Li, Y. (2019). What can resnet learn efficiently, going beyond kernels? *Advances*
403 *in Neural Information Processing Systems*, **32**.
- 404 Anandkumar, A., Deng, Y., Ge, R. and Mobahi, H. (2017). Homotopy analysis for tensor pca. In
405 *Conference on Learning Theory*. PMLR.
- 406 Arnaboldi, L., Dandi, Y., Krzakala, F., Pesce, L. and Stephan, L. (2024). Repetita iuvant: Data repe-
407 tition allows sgd to learn high-dimensional multi-index functions. In *High-dimensional Learning*
408 *Dynamics 2024: The Emergence of Structure and Reasoning*.
- 409 Arous, G. B., Gheissari, R. and Jagannath, A. (2021). Online stochastic gradient descent on non-
410 convex losses from high-dimensional inference. *Journal of Machine Learning Research*, **22**
411 1–51.
- 412 Arous, G. B., Wein, A. S. and Zadik, I. (2020). Free energy wells and overlap gap property in sparse
413 pca. In *Conference on Learning Theory*. PMLR.
- 414 Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z. and Wu, D. (2023). Learning in the presence of low-
415 dimensional structure: a spiked random matrix perspective. *Advances in Neural Information*
416 *Processing Systems*, **36**.
- 417 Bach, F. (2017). Breaking the curse of dimensionality with convex neural networks. *Journal of*
418 *Machine Learning Research*, **18** 1–53.
- 419 Baik, J., Arous, G. B. and P ech e, S. (2005). Phase transition of the largest eigenvalue for nonnull
420 complex sample covariance matrices. *The Annals of Probability*, **33** 1643 – 1697.
- 421 Bandeira, A. S., El Alaoui, A., Hopkins, S., Schramm, T., Wein, A. S. and Zadik, I. (2022). The
422 franz-parisi criterion and computational trade-offs in high dimensional statistics. *Advances in*
423 *Neural Information Processing Systems*, **35** 33831–33844.
- 424 Bietti, A., Bruna, J., Sanford, C. and Song, M. J. (2022). Learning single-index models with shallow
425 neural networks. *Advances in Neural Information Processing Systems*, **35** 9768–9783.
- 426 Biroli, G., Cammarota, C. and Ricci-Tersenghi, F. (2020). How to iron out rough landscapes and get
427 optimal performances: averaged gradient descent and its application to tensor pca. *Journal of*
428 *Physics A: Mathematical and Theoretical*, **53** 174003.
- 429 Cand es, E. J., Romberg, J. and Tao, T. (2006). Robust uncertainty principles: Exact signal reconstruc-
430 tion from highly incomplete frequency information. *IEEE Transactions on information theory*, **52**
431 489–509.
- 432 Chen, M., Bai, Y., Lee, J. D., Zhao, T., Wang, H., Xiong, C. and Socher, R. (2020). Towards under-
433 standing hierarchical learning: Benefits of neural representations. *Advances in Neural Information*
434 *Processing Systems*, **33** 22134–22145.
- 435 Chen, S. and Meka, R. (2020). Learning polynomials in few relevant dimensions. In *Conference on*
436 *Learning Theory*. PMLR.
- 437 Chizat, L., Oyallon, E. and Bach, F. (2019). On lazy training in differentiable programming. *Advances*
438 *in neural information processing systems*, **32**.
- 439 Choo, D. and d’Orsi, T. (2021). The complexity of sparse tensor pca. *Advances in Neural Information*
440 *Processing Systems*, **34** 7993–8005.
- 441 Damian, A., Lee, J. and Soltanolkotabi, M. (2022). Neural networks can learn representations with
442 gradient descent. In *Conference on Learning Theory*. PMLR.

- 443 Damian, A., Nichani, E., Ge, R. and Lee, J. D. (2023). Smoothing the landscape boosts the signal for
444 sgd: Optimal sample complexity for learning single index models. *Advances in Neural Information*
445 *Processing Systems*, **36**.
- 446 Damian, A., Pillaud-Vivien, L., Lee, J. D. and Bruna, J. (2024). The computational complexity of
447 learning gaussian single-index models. *arXiv preprint arXiv:2403.05529*.
- 448 Dandi, Y., Krzakala, F., Loureiro, B., Pesce, L. and Stephan, L. (2023). How two-layer neural net-
449 works learn, one (giant) step at a time. In *NeurIPS 2023 Workshop on Mathematics of Modern*
450 *Machine Learning*.
- 451 Dandi, Y., Troiani, E., Arnaboldi, L., Pesce, L., Zdeborova, L. and Krzakala, F. (2024). The benefits
452 of reusing batches for gradient descent in two-layer networks: Breaking the curse of information
453 and leap exponents. In *Forty-first International Conference on Machine Learning*.
- 454 d’Aspremont, A., Ghaoui, L., Jordan, M. and Lanckriet, G. (2004). A direct formulation for sparse
455 pca using semidefinite programming. *Advances in neural information processing systems*, **17**.
- 456 Donoho, D. L., Maleki, A. and Montanari, A. (2009). Message-passing algorithms for compressed
457 sensing. *Proceedings of the National Academy of Sciences*, **106** 18914–18919.
- 458 Dudeja, R. and Hsu, D. (2018). Learning single-index models in gaussian space. In *Conference On*
459 *Learning Theory*. PMLR.
- 460 Dudeja, R. and Hsu, D. (2021). Statistical query lower bounds for tensor pca. *Journal of Machine*
461 *Learning Research*, **22** 1–51.
- 462 Fan, J., Yang, Z. and Yu, M. (2023). Understanding implicit regularization in over-parameterized
463 single index model. *Journal of the American Statistical Association*, **118** 2315–2328.
- 464 Feldman, V. (2017). A general characterization of the statistical query complexity. In *Conference on*
465 *learning theory*. PMLR.
- 466 Feldman, V., Grigorescu, E., Reyzin, L., Vempala, S. S. and Xiao, Y. (2017). Statistical algorithms
467 and a lower bound for detecting planted cliques. *Journal of the ACM (JACM)*, **64** 1–37.
- 468 Folland, G. B. (2001). How to integrate a polynomial over a sphere. *The American Mathematical*
469 *Monthly*, **108** 446–448.
- 470 Gamarnik, D. and Zadik, I. (2017). Sparse high-dimensional linear regression. algorithmic barriers
471 and a local search algorithm. *arXiv preprint arXiv:1711.04952*.
- 472 Ghorbani, B., Mei, S., Misiakiewicz, T. and Montanari, A. (2019). Limitations of lazy training of
473 two-layers neural network. *Advances in Neural Information Processing Systems*, **32**.
- 474 Girshick, R., Donahue, J., Darrell, T. and Malik, J. (2014). Rich feature hierarchies for accurate
475 object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer*
476 *vision and pattern recognition*.
- 477 Härdle, W., Müller, M., Sperlich, S., Werwatz, A. et al. (2004). *Nonparametric and semiparametric*
478 *models*, vol. 1. Springer.
- 479 Hristache, M., Juditsky, A. and Spokoiny, V. (2001). Direct estimation of the index coefficient in a
480 single-index model. *Annals of Statistics* 595–623.
- 481 Ichimura, H. (1993). Semiparametric least squares (sls) and weighted sls estimation of single-index
482 models. *Journal of econometrics*, **58** 71–120.
- 483 Jacot, A., Gabriel, F. and Hongler, C. (2018). Neural tangent kernel: Convergence and generalization
484 in neural networks. *Advances in neural information processing systems*, **31**.
- 485 Jin, C., Ge, R., Netrapalli, P., Kakade, S. M. and Jordan, M. I. (2017). How to escape saddle points
486 efficiently. In *International conference on machine learning*. PMLR.

- 487 Johnstone, I. M. and Lu, A. Y. (2009). On consistency and sparsity for principal components analysis
488 in high dimensions. *Journal of the American Statistical Association*, **104** 682–693.
- 489 Joshi, N., Misiakiewicz, T. and Srebro, N. (2024). On the complexity of learning sparse functions
490 with statistical and gradient queries. *arXiv preprint arXiv:2407.05622*.
- 491 Klenke, A. and Mattner, L. (2010). Stochastic ordering of classical discrete distributions. *Advances
492 in Applied probability*, **42** 392–410.
- 493 Lee, J. D., Oko, K., Suzuki, T. and Wu, D. (2024). Neural network learns low-dimensional polyno-
494 mials with sgd near the information-theoretic limit. *arXiv preprint arXiv:2406.01581*.
- 495 Lu, Y. M. and Li, G. (2020). Phase transitions of spectral initialization for high-dimensional non-
496 convex estimation. *Information and Inference: A Journal of the IMA*, **9** 507–541.
- 497 MacCullagh, P. and Nelder, J. (1989). Generalized linear models. *Monographs on statistics and
498 applied probability (37)*.
- 499 Mondelli, M. and Montanari, A. (2018). Fundamental limits of weak recovery with applications to
500 phase retrieval. In *Conference On Learning Theory*. PMLR.
- 501 Neykov, M., Wang, Z. and Liu, H. (2016). Agnostic estimation for misspecified phase retrieval
502 models. *Advances in Neural Information Processing Systems*, **29**.
- 503 O’Donnell, R. (2014). *Analysis of boolean functions*. Cambridge University Press.
- 504 Raskutti, G., J Wainwright, M. and Yu, B. (2012). Minimax-optimal rates for sparse additive models
505 over kernel classes via convex programming. *Journal of machine learning research*, **13**.
- 506 Refinetti, M., Goldt, S., Krzakala, F. and Zdeborová, L. (2021). Classifying high-dimensional gaus-
507 sian mixtures: Where kernel methods fail and neural networks succeed. In *International Conference
508 on Machine Learning*. PMLR.
- 509 Szekli, R. (2012). *Stochastic ordering and dependence in applied probability*, vol. 97. Springer
510 Science & Business Media.
- 511 Vaskevicius, T., Kanade, V. and Rebeschini, P. (2019). Implicit regularization for optimal sparse
512 recovery. *Advances in Neural Information Processing Systems*, **32**.
- 513 Wang, T., Berthet, Q. and Samworth, R. J. (2016). Statistical and computational trade-offs in estima-
514 tion of sparse principal components.
- 515 Zhao, P., Yang, Y. and He, Q.-C. (2022). High-dimensional linear regression via implicit regulariza-
516 tion. *Biometrika*, **109** 1033–1046.

517	Contents	
518	1 Introduction	1
519	2 Problem Setup	4
520	3 Overview of techniques	4
521	4 Gradient-based Algorithm for Uniform Prior	6
522	4.1 Gradient-based Training Algorithm (Algorithm 1)	6
523	4.2 Feature Alignment and Statistical Complexity	7
524	4.2.1 Online SGD with Batch Reusing	8
525	4.2.2 Label Transformation via Modified Loss	8
526	5 Exploiting the Structure: Algorithm for Sparse Prior	8
527	5.1 Algorithm design: How to leverage sparsity?	9
528	5.2 Sample Complexity Analysis for Sparse Prior	9
529	A Numerical Experiments	15
530	B Notation and Preliminaries	17
531	B.1 Background on Hermite Polynomials	18
532	C Supplementary Proofs for the Main Context	18
533	C.1 Proofs for Section 3	18
534	C.2 Proofs for Examples of Oracle Function	19
535	C.2.1 Batch-reusing for polynomial link function	19
536	C.2.2 Modified loss for general $s^* \geq 1$	22
537	D Proof Sketch of the Main Theorem for Uniform Prior	24
538	E Proof of the Main Theorem for the Uniform Prior	25
539	E.1 Properties of the Gradient Step	26
540	E.2 Proof of the Main Theorem for Uniform Prior	27
541	E.3 Proof of Key Results	32
542	F Proof of the Main Theorem for the Sparse Prior	38
543	F.1 Proof Outline and Preliminaries	38
544	F.2 Properties of the Gradient Step	40
545	F.3 Proof of the Main Theorem	41
546	F.4 Proof of the Key Results	46
547	F.5 Proofs for Technical Results in the Sparse Case	51
548	G Statistical Query Lower Bound for Sparse Signal Recovery	55

549	H Supporting Lemmas on Moment Calculations	59
550	I Technical Results	64
551	I.1 Technical Results for Hermite Tensor	64
552	I.2 Technical Results for Uniform Distribution on the Sphere	69
553	I.3 Technical Results on Polarized Random Vectors	71

554	J Auxiliary Lemmas	74
-----	---------------------------	-----------

555 **A Numerical Experiments**

556 We conduct extensive simulation experiments to validate the sample complexity result established in
 557 Theorem 4.2. In specific, for a fixed Gaussian single-index model, we run Algorithm 1 extensively
 558 over a variety of problem instances with diverse scales and report the average accuracy in terms of
 559 the alignment. We lay out the details of the experiment setting as follows.

560 • **Gaussian single-index model.** We focus on the Gaussian single-index model introduced in (2.1)
 561 with a deterministic link function p and Gaussian additive noise. Here we set $p(x) = x^2 \cdot \exp(-x^2) + \epsilon$
 562 where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ with $\sigma^2 = 0.5$. As shown in Example 2.2, the generative exponent of the function
 563 p is $s^*(p) = 4$. In addition, the signal parameter θ^* is uniformly sampled from the unit sphere in \mathbb{R}^d .

564 • **Neural network architecture.** We adopt the two-layer neural network introduced in Section 2 with
 565 M set to 15 and $a_m = 1$ for all $m \in [M]$ in all experiments. Since $s^*(p) = 4$, we set the activation
 566 function as $\sigma(x) = h_4(x)$, i.e., the fourth-order Hermite polynomial.

567 • **Training using Algorithm 1.** In Algorithm 1, we set $\psi(x, y) = y \cdot \sigma'(x)$, as stated in Example 4.6.
 568 Such a ψ is justified by considering the following alignment loss:

$$L(\theta) = 1 - y \cdot f(z; \theta, \mathbf{a}) = 1 - \sum_{m=1}^M a \cdot \sigma(\langle z, \theta_m \rangle) \cdot y, \quad (\text{A.1})$$

569 where recall that each entry of \mathbf{a} is equal to a . As a result, by (A.1) we have

$$a^{-1} \cdot \nabla_{\theta_m} L(\theta) = y \cdot \sigma'(\langle z, \theta_m \rangle) \cdot z = \psi(\langle z, \theta_m \rangle, y) \cdot z.$$

570 As a result, we can alternatively interpret the gradients in Algorithm 1 as those with respect to the
 571 alignment loss $L(\theta)$. Thus, the choice of a does not matter in this case, and we set $a = 1$ for simplicity.
 572 Furthermore, other details of Algorithm 1 are specified as follows:

- 573 • The parameters $\{\theta_n\}_{m \in [M]}$ are initialized as i.i.d. random vectors in \mathbb{R}^d uniformly sampled
 574 from the unit sphere.
- 575 • We fix $M = 15$, $a = 1$, $\eta = 3$, $T = 24$, and $L = 500$ throughout all experiments with
 576 different values of n and d .
- 577 • We enumerate n and d over a grid with $d \in [32, 499]$ and $n \in [5 \times 10^3, 3 \times 10^6]$. Note that
 578 $\log d \in (3, 7)$, our choice of T satisfies the requirement in Theorem 4.2.

579 • **Choices of (d, n) .** We select 40 different values of d and 30 different values of n within the ranges
 580 $d \in [32, 499]$ and $n \in [5 \times 10^3, 3 \times 10^6]$, respectively. These values form an evenly spaced grid in
 581 terms of $\log n$ and $\log d$. See Figure 2 for an illustration.

582 • **Evaluation.** We report the accuracy of Algorithm 1 based on 25 repeated experiments for every
 583 choice of (n, d) . We report two types of accuracy metrics:

- 584 (i) Average accuracy: We report $M^{-1} \sum_{m=1}^M |\langle \theta_m, \theta^* \rangle|$ in each experiment and then average
 585 over the 25 experiments.

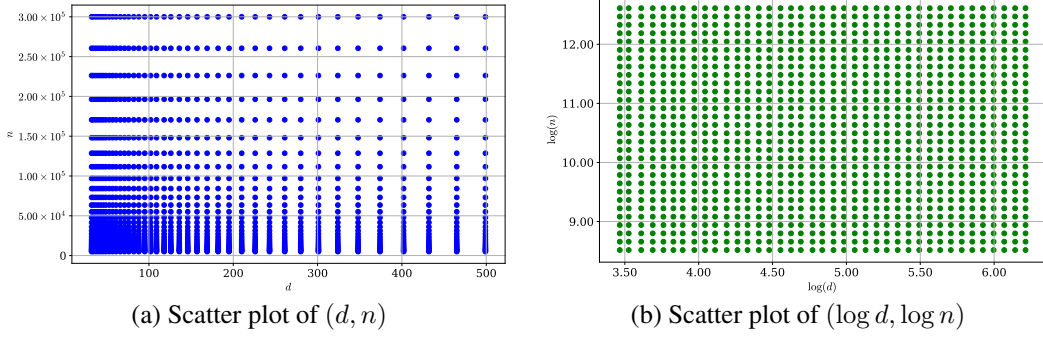


Figure 2: Scatter plots of (d, n) and $(\log d, \log n)$. In (a) we plot n against d and in (b) we plot $\log n$ against $\log d$. As shown in (b), we choose n and d such that they form an evenly-spaced grid after logarithm.

586 (ii) Top-8 accuracy: Given $\{\theta_m\}_{m \in [M]}$ returned by the algorithm, we sort the alignment values
 587 $\{|\langle \theta_m, \theta^* \rangle|\}_{m \in [M]}$. Then we report the average of the largest 8 numbers. The rationale is
 588 that if the top-8 accuracy is close to one, at least half of the neurons correctly find θ^* .

589 **Contour plots.** After calculating these two versions of accuracy for every (d, n) pair, we generate
 590 the contour plots based on $(\log d, \log n, \text{acc}(d, n))$, where $\text{acc}(d, n)$ is one of the two versions of
 591 average accuracy introduced above. We report these two contour plots in Figure 3 and Figure 4, where
 592 in Figure 3 we zoom in to a smaller range of d for better visualization. In these plots, points with the
 593 same color indicate $(\log d, \log n)$ with the same level of accuracy.

594 **Validate $\tilde{\Theta}(d^{s^*/2})$ sample complexity.** As shown in these figures, the average accuracy and the top-8
 595 accuracy clearly exhibit a **linear relationship**. That is, for a fixed accuracy level δ , (d, n) satisfying
 596 $\text{acc}(d, n) = \delta$ is a line segment. That is, $\log n = c_1 \cdot \log d + c_2$. To determine c_1 and c_2 , we further
 597 fit linear models for $(\log d, \log n)$ with the same accuracy level δ , where $\delta \in \{0.6, 0.7, 0.8\}$. For both
 598 the average accuracy and the top-8 accuracy, the coefficient c_1 in the linear models is close to 2. We
 599 report the linear models corresponding to different accuracy levels in Table 1. This finding indicates
 600 that $n \propto d^2$. Note that $s^* = 4$. Moreover, since we compute the accuracy for all (d, n) on the grid.
 601 The fact that $c_1 \approx 2$ indicates that the $\tilde{\Theta}(d^{s^*/2})$ sample complexity is sharp.

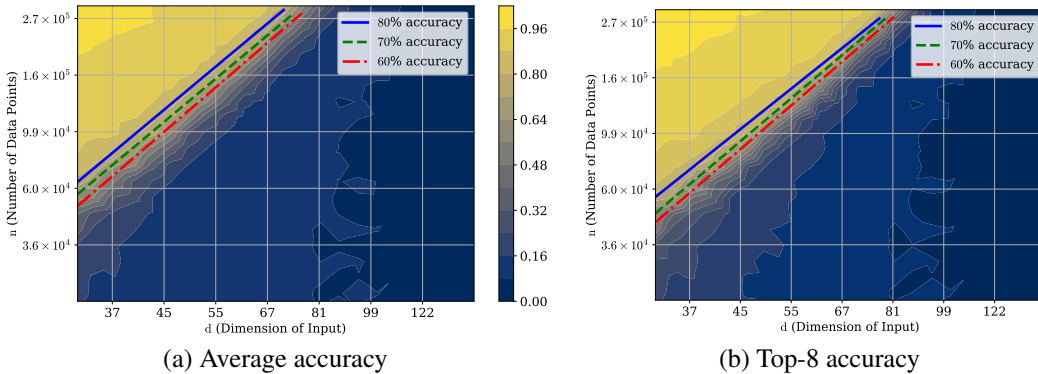


Figure 3: The contour plots of $(\log d, \log n, \text{acc}(d, n))$, where $\text{acc}(d, n)$ is either the average accuracy and top-8 accuracy. Here we zoom in to a smaller subset of d 's for better visualization. We also plot the lines containing $(\log d, \log n)$ with the same accuracy level among $\{0.6, 0.7, 0.8\}$. The slopes of these lines are all close to 2. This indicates that $n \approx d^2$ samples are sufficient and necessary for accurate estimation.

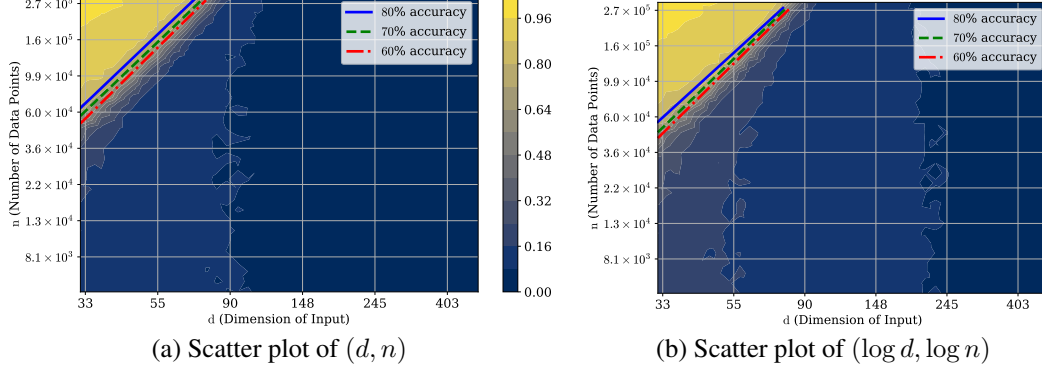


Figure 4: The contour plots of $(\log d, \log n, \text{acc}(d, n))$, where $\text{acc}(d, n)$ is either the average accuracy and top-8 accuracy. We also plot the lines containing $(\log d, \log n)$ with the same accuracy level among $\{0.6, 0.7, 0.8\}$. The slopes of these lines are all close to 2. This indicates that $n \approx d^2$ samples are sufficient and necessary for accurate estimation.

Table 1: Fitted linear equations of the form $\log n = c_1 \cdot \log d + c_2$ for n, d with the desired accuracy level. Notably, the slopes of these equations are all close to $s^*/2 = 2$, which shows that $n \propto d^{s^*/2}$.

Accuracy level	Average accuracy	Top-8 accuracy
0.8	$\log(n) = 1.9058 \cdot \log(d) + 4.4516$	$\log(n) = 1.8201 \cdot \log(d) + 4.6218$
0.7	$\log(n) = 1.9103 \cdot \log(d) + 4.3273$	$\log(n) = 1.9343 \cdot \log(d) + 4.0790$
0.6	$\log(n) = 1.9640 \cdot \log(d) + 4.0361$	$\log(n) = 1.9653 \cdot \log(d) + 3.8901$

602 B Notation and Preliminaries

603 **Notations.** We use \mathbb{N} to denote the set of positive integers and \mathbb{N}_0 to denote the set of nonnegative
604 integers. For vector $z \in \mathbb{R}^d$, we denote by $\mathbb{R}_n[z]$ the set of polynomials of degree at most n in z with
605 real coefficients. For $s \in \mathbb{N}$, we denote by Π_s the symmetric group of all permutations of $[s]$. We
606 denote by $\mathcal{N}_d(\cdot)$ and $\mathcal{N}(\cdot)$ the standard normal distribution in \mathbb{R}^d and \mathbb{R} , respectively.

607 For two tensors $S \in (\mathbb{R}^d)^{\otimes s}$ and $T \in (\mathbb{R}^d)^{\otimes t}$ where $s \geq t$,

$$(S[T])_{j_1, \dots, j_{s-t}} := \sum_{i_1, \dots, i_t=1}^d S_{j_1, \dots, j_{s-t}, i_1, \dots, i_t} T_{i_1, \dots, i_t}.$$

608 Here, $S[T]$ produces a tensor of order $s - t$ and dimension d . We also define the symmetrization
609 operation for a tensor $T \in (\mathbb{R}^d)^{\otimes t}$ as

$$\text{Sym}(T)_{i_1, \dots, i_t} := \frac{1}{t!} \sum_{\pi \in \Pi_t} T_{i_{\pi(1)}, \dots, i_{\pi(t)}}.$$

610 The followings are some notations for the relationship between two quantities (or matrices):

611 $a \simeq b$: There exists a constant $C = O(1)$ such that $a \leq Cb$ and $b \leq Ca$. Note that a and b should
612 have the same sign. $a = \Theta(b)$ also has the same meaning.

613 $a \approx b$: $a \leq \text{polylog}(d) \cdot b$ and $b \leq \text{polylog}(d) \cdot a$, and the same for $a = \tilde{\Theta}(b)$.

614 $a \lesssim b$: There exists a constant $C = O(1)$ such that $a \leq Cb$, and the same for $a = \Omega(b)$. The use of
615 $a \gtrsim b$ is similar.

616 $a \lesssim b$: $a \leq \text{polylog}(d) \cdot b$, and the same for $a = \tilde{\Omega}(b)$. The use of $a \gtrsim b$ and $a = \tilde{\Omega}(b)$ is similar.

617 $a \ll b$: $a \leq (\text{polylog}(d))^{-1} \cdot b$. The use of $a \gg b$ is similar.

618 In addition, we denote by $a = b \pm \varepsilon$, $a \simeq b \pm \varepsilon$, $a \approx b \pm \varepsilon$ that $b - \varepsilon \leq a \leq b + \varepsilon$, $a - \varepsilon \lesssim b \lesssim a + \varepsilon$,
619 $a - \varepsilon \lesssim b \lesssim a + \varepsilon$, respectively.

620 For square matrices A and B , $A \preceq B$ means that $B - A$ is positive semi-definite, and $A \lesssim B$ means
621 that there exists a constant $C = O(1)$ such that $C \cdot B - A$ is positive semidefinite.

622 **B.1 Background on Hermite Polynomials**

623 The probabilist’s Hermite polynomials satisfy the following recurrence relations

$$h_s(x)' = \sqrt{s} \cdot h_{s-1}(x), \quad x \cdot h_s(x) = \sqrt{s+1} \cdot h_{s+1}(x) + \sqrt{s} \cdot h_{s-1}(x), \quad (\text{B.1})$$

624 where we adopt the convention that $h_{-1}(x) \equiv 0$.

625 For any function $f \in L^2(\mathcal{N}(0, 1))$, its Hermite expansion is given by

$$f(x) = \sum_{s=0}^{\infty} \hat{f}_s \cdot h_s(x),$$

626 where we denote by \hat{f}_s the s -th coefficient of the Hermite expansion of f .

627 **Gaussian noise operator.** For $\rho \in [-1, 1]$, define the Gaussian noise operator as

$$U_\rho f(x) = \mathbb{E}_{x' \sim \mathcal{N}(0,1)} [f(\rho x + \sqrt{1-\rho^2} \cdot x')].$$

628 Proposition 11.37 of O’Donnell (2014) shows that the Hermite expansion of $U_\rho f$ is given by

$$U_\rho f(x) \stackrel{L^2(\mathcal{N}(0,1))}{=} \sum_{s=0}^{\infty} \rho^s \cdot \hat{f}_s \cdot h_s(x).$$

629 A direct implication of this identity is

$$\mathbb{E}_{x \sim \mathcal{N}(0,1)} [U_\rho f(x) g(x)] = \mathbb{E}_{x \sim \mathcal{N}(0,1)} [f(x) U_\rho g(x)] = \sum_{s=0}^{\infty} \rho^s \hat{f}_s \hat{g}_s. \quad (\text{B.2})$$

630 As a result, for any fixed $w, \theta \in \mathbb{S}^{d-1}$, it holds that

$$\mathbb{E}_{z \sim \mathcal{N}(0, I_d)} [f(\langle w, z \rangle) h_s(\langle \theta, z \rangle)] = \mathbb{E}_{x \sim \mathcal{N}(0,1)} [U_\rho f \cdot h_s(x)] = \langle w, \theta \rangle^s \cdot \hat{f}_s. \quad (\text{B.3})$$

631 **Hermite tensor.** Corresponding to the Hermite polynomials defined for scalar variables, we define
632 the Hermite tensors over $z \in \mathbb{R}^d$:

$$\mathbf{h}_s(z) := \frac{(-1)^s}{\sqrt{s!}} \cdot e^{\|z\|_2^2/2} \cdot \nabla^s e^{-\|z\|_2^2/2} \in (\mathbb{R}^d)^{\otimes s}, \text{ for } s \geq 0.$$

633 The scalar-valued Hermite polynomials and the tensor-valued Hermite tensors are related as follows:

$$h_s(\langle \theta, z \rangle) = \mathbf{h}_s(z) [\theta^{\otimes s}], \quad \forall \theta \in \mathbb{S}^{d-1}. \quad (\text{B.4})$$

634 Now let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a s -times differentiable function such that for all $k \leq s$, every component of
635 $\nabla^k f$ belongs to $L^2(\mathcal{N}(0, I_d))$. Then it follows from integration by parts that

$$\mathbb{E}_{z \sim \mathcal{N}(0, I_d)} [f(z) \mathbf{h}_s(z)] = \frac{1}{\sqrt{s!}} \cdot \mathbb{E}_{z \sim \mathcal{N}(0, I_d)} [\nabla^s f(z)]. \quad (\text{B.5})$$

636 This is a version of Stein’s lemma for tensor-valued functions.

637 **C Supplementary Proofs for the Main Context**

638 **C.1 Proofs for Section 3**

639 In this section, we first argue why $\mathbb{E}_{x \sim \mathcal{N}(0,1)} [y \zeta_{s^*}(y)] = 0$ is the major difficulty for vanilla (stochas-
640 tic) gradient descent to achieve the information-theoretical lower bound $O(d)$ (the same for SQ lower
641 bound) when the information exponent q^* is larger than 2. It has been shown by Damian et al. (2024)
642 that the generative exponent s^* for polynomial model is either 1 or 2. Consider the information
643 exponent $q^* > 2$. We have the following lemma saying that the correlation $\mathbb{E}_{x \sim \mathcal{N}(0,1)} [y \zeta_s(y)] = 0$
644 for any $s < q^*$.

645 **Lemma C.1.** Recall that ζ_s is the coefficient function for degree s in the decomposition of the
646 likelihood ratio $\mathbb{P}(x, y)/\mathbb{Q}(x, y)$ in (2.2). For any $q^* \geq 2$, consider the Gaussian single-index model
647 given by $y = \beta_0 + \sum_{p \geq p^*} \beta_p h_p(x)$ with $x \sim \mathcal{N}(0, 1)$. Then for any $1 \leq s < p^*$, $\mathbb{E}_{\mathbb{Q}}[y \zeta_s(y)] = 0$.

648 *Proof.* The proof can be done by noting that $\zeta_s(y) = \mathbb{E}_{\mathbb{Q}}[\mathbb{P}(x, y)/\mathbb{Q}(x, y) \cdot h_s(x) \mid y]$, and

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}}[y \zeta_s(y)] &= \mathbb{E}_{\mathbb{Q}} \left[y \cdot \mathbb{E}_{\mathbb{Q}} \left[\frac{\mathbb{P}(x, y)}{\mathbb{Q}(x, y)} \cdot h_s(x) \mid y \right] \right] = \mathbb{E}_{\mathbb{P}} [y \cdot h_s(x)] \\ &= \beta_0 \cdot \mathbb{E}_{x \sim \mathcal{N}(0, 1)} [h_s(x)] + \sum_{i \geq p^*} \beta_i \cdot \mathbb{E}_{x \sim \mathcal{N}(0, 1)} [h_i(x) h_s(x)] = 0, \end{aligned}$$

649 where the second equality follows from the independence between x and y under \mathbb{Q} . \square

650 Therefore, the first nonzero term in the informative queries of (3.1) is of order at least p^* . This gives
651 rise to sample complexity $d^{p^* - 1}$ for vanilla online SGD (Arous et al., 2021) and $d^{p^*/2}$ for SGD after
652 smoothing the landscape (Damian et al., 2023). This sample complexity $d^{p^*/2}$ matches the correlated
653 statistical query (CSQ) lower bound with gradient of the form $y\phi(z)$ (Abbe et al., 2023; Damian
654 et al., 2022).

655 C.2 Proofs for Examples of Oracle Function

656 Here we complete the discussions of the specific examples of ψ in Example 4.4 and Example 4.6.

657 C.2.1 Batch-reusing for polynomial link function

658 We consider a polynomial link function $y = p(x) = \sum_{q^* \leq q' \leq q} \beta_{q'} h_{q'}(x)$ for general $q^* \in \mathbb{N}$ and
659 $\beta_{q'} \in \mathbb{R}$, where q^* is the information exponent of the link function, and we also denote it by $q^*(p)$
660 in the sequel. For batch-reusing, we take $\psi(y, x) = y\sigma'(x) + y\sigma'(x + y\sigma'(x))$, where activation
661 function $\sigma(x)$ satisfies that

$$\sigma(x) = \sum_{j=0}^{C_q} \alpha_j \cdot h_j(x), \quad \sigma'(x) = \sum_{j=1}^{C_q} \sqrt{j} \cdot \alpha_j \cdot h_{j-1}(x). \quad (\text{C.1})$$

662 Here the degree $C_q \in \mathbb{N}_+$ only depends on the degree q of the link function and is specified later,
663 and each coefficient $\alpha_j \sim \text{Unif}([0, 1])$. The second equality in (C.1) follows from the property
664 of Hermite polynomials in (B.1). The error term $\text{err}_{m,l,i}^{(t)}$ now comes from the difference between
665 $\psi(y_i^{(t)}, \langle w_{m,l}^{(t)}, z_i^{(t)} \rangle) \cdot z_i^{(t)}$ and the exact form of the update step obtained from two consecutive
666 gradient descent steps on the same data under the square loss. More specifically, let us consider a
667 single neuron whose weight is $w_{m,l}$ and a single data point (z_i, y_i) . Here we omit the time index t
668 for convenience. Then two gradient descent step on (z_i, y_i) gives

$$\begin{aligned} -g_{m,l}^{\text{Re}}(z_i, y_i) &= (y_i - f(z_i; \{w_{m,l}\}_{m \in [M]})) \cdot \sigma'(\langle w_{m,l}, z_i \rangle) \cdot z_i \\ &\quad + (y_i - f(z_i; \{w_{m,l}^+\}_{m \in [M]})) \cdot \sigma'(\langle w_{m,l}^+, z_i \rangle) \cdot z_i, \end{aligned} \quad (\text{C.2})$$

669 where $w_{m,l}^+ = w_{m,l} + \eta_i^{\text{Re}} \cdot (y_i - f(z_i; \{w_{m,l}\}_{m \in [M]})) \cdot \sigma'(\langle w_{m,l}, z_i \rangle) \cdot z_i$. Here η_i^{Re} is the learning
670 rate for batch reusing, different from the learning rate η in our algorithm. More specifically, to
671 fit the gradient form (C.2) into our general framework with oracle function $\psi(y, x)$, we take the
672 batch-reusing learning rate $\eta_i^{\text{Re}} = 1/\|z_i\|_2^2$. Then the error term is given by

$$\text{err}_{m,l,i} = -g_{m,l}^{\text{Re}}(z_i, y_i) - \psi(y_i, \langle w_{m,l}, z_i \rangle) = \text{err}_{m,l,i,1} + \text{err}_{m,l,i,2} + \text{err}_{m,l,i,3}, \quad (\text{C.3})$$

673 where $\text{err}_{m,l,i,1}$, $\text{err}_{m,l,i,2}$, and $\text{err}_{m,l,i,3}$ are given by

$$\begin{aligned} \text{err}_{m,l,i,1} &= -f(z_i; \{w_{m,l}\}_{m \in [M]}) \cdot \sigma'(\langle w_{m,l}, z_i \rangle), \\ \text{err}_{m,l,i,2} &= -f(z_i; \{w_{m,l}^+\}_{m \in [M]}) \cdot \sigma'(\langle w_{m,l}^+, z_i \rangle), \\ \text{err}_{m,l,i,3} &= y_i \sigma'(\langle w_{m,l}, z_i \rangle) + y_i \sigma'(\langle w_{m,l}, z_i \rangle) + \text{err}_{m,l,i,1} \\ &\quad - y_i \sigma'(\langle w_{m,l}, z_i \rangle + y_i \sigma'(\langle w_{m,l}, z_i \rangle)). \end{aligned} \quad (\text{C.4})$$

674 *Proof of Corollary 4.5.* To prove Corollary 4.5, it suffices to show that (i) Assumption 4.1 holds, and
 675 (ii) the event \mathcal{E} holds with the desired high probability. In the following, we first verify Assumption 4.1,
 676 and then check the event \mathcal{E} .

677 **Verifying Assumption 4.1.** Note that the fact of $y = p(x)$ being a polynomial immediately implies
 678 that both the square-integrable condition (Assumption 4.1(a)) and the polynomial-like tail condition
 679 (Assumption 4.1(c)) are satisfied. It remains to check the high-pass condition (Assumption 4.1(b)).
 680 Since now $s^* \leq 2$, we only need to check the condition that $|\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)]| > 0$.

681 **Case 1:** $s^* = 2$. In this case, we have that

$$\begin{aligned} \hat{\psi}_1(y) &= \mathbb{E}_{x \sim \mathcal{N}} \left[x \cdot \left(y\sigma'(x) + y\sigma'(x + y\sigma'(x)) \right) \right] \\ &= y \cdot \mathbb{E}_{x \sim \mathcal{N}} \left[\sum_{j=1}^{C_q} \sqrt{j} \cdot \alpha_j \cdot x \cdot h_{j-1}(x) + \sum_{j=1}^{C_q} \sqrt{j} \cdot \alpha_j \cdot x \cdot h_{j-1}(x + y\sigma'(x)) \right] \end{aligned} \quad (\text{C.5})$$

682 For the first summation in (C.5), only the first summand is nonzero, so we obtain

$$y \cdot \mathbb{E}_{x \sim \mathcal{N}} \left[\sum_{j=1}^{C_q} \sqrt{j} \cdot \alpha_j \cdot x \cdot h_{j-1}(x) \right] = \sqrt{2}\alpha_2 \cdot y. \quad (\text{C.6})$$

683 For the second summation in (C.5), we have the following expansion,

$$\begin{aligned} &y \cdot \mathbb{E}_{x \sim \mathcal{N}} \left[\sum_{j=1}^{C_q} j \cdot \alpha_j \cdot x \cdot h_{j-1}(x + y\sigma'(x)) \right] \\ &= y \cdot \mathbb{E}_{x \sim \mathcal{N}} \left[\sum_{j=1}^{C_q} j \cdot \alpha_j \cdot x \cdot \sum_{k=0}^{j-1} r_{j-1,k} \cdot h_{j-k-1}(x) \cdot (y\sigma'(x))^k \right] \\ &= \sum_{k=0}^{C_q-1} \underbrace{\left\{ \sum_{j=k+1}^{C_q} j \cdot \alpha_j \cdot r_{j-1,k} \cdot \mathbb{E}_{x \sim \mathcal{N}} \left[x \cdot h_{j-k-1}(x) \cdot (\sigma'(x))^k \right] \right\}}_{:= \varsigma_k(\alpha)} \cdot y^{k+1}. \end{aligned} \quad (\text{C.7})$$

684 where $\varsigma_0(\alpha), \dots, \varsigma_{C_q-1}(\alpha)$ are just polynomials of $\alpha = (\alpha_1, \dots, \alpha_{C_q})$ (recall the definition of $\sigma'(x)$
 685 in (C.1)) and each $r_{j-1,k}$ is a positive number. Combining (C.6) and (C.7), we get the following
 686 decomposition of $\hat{\psi}_1(y)$:

$$\hat{\psi}_1(y) = \sqrt{2}\alpha_2 \cdot y + \sum_{k=0}^{C_q-1} \varsigma_k(\alpha) \cdot y^{k+1}.$$

687 Further using $y = p(x)$ and the definition of $\zeta_2(y)$, we get

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}}[\zeta_2(y) \cdot \hat{\psi}_1(y)] &= \mathbb{E}_{\mathbb{P}}[h_2(x) \cdot \hat{\psi}_1(y)] \\ &= \sqrt{2}\alpha_2 \cdot \mathbb{E}_{x \sim \mathcal{N}} [h_2(x) \cdot p(x)] + \sum_{k=0}^{C_q-1} \varsigma_k(\alpha) \cdot \mathbb{E}_{x \sim \mathcal{N}} [h_2(x) \cdot p(x)^{k+1}]. \end{aligned}$$

688 According to Proposition 5 of Lee et al. (2024), we can set $C_q \in \mathbb{N}^+$ (only depending on q) such
 689 that there exists a smallest $I \leq C_q$ such that the information exponent $q^*(p^I) \leq 2$. We notice that
 690 in this case $s^*(p) = 2$, where we abuse the notation and let $s^*(p)$ be the generative exponent of the
 691 polynomial p . In fact, $s^*(p) = 1$ means $\mathbb{E}_{\mathbb{P}}[\mathcal{T}(y) \cdot h_1(x)] \equiv 0$ for all label transformation \mathcal{T} . Hence,
 692 the only possibility is that $q^*(p^I) = 2$ since p^I is just a special case of label transformation and we
 693 cannot get any first-order term from p^I . Therefore, we further simplify the target quantity as

$$\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] = \sqrt{2}\alpha_2 \cdot \underbrace{\mathbb{E}_{x \sim \mathcal{N}} [h_2(x) \cdot p(x)]}_{:= b_1} + \sum_{k=I}^{C_q} \underbrace{\varsigma_{k-1}(\alpha) \cdot \mathbb{E}_{x \sim \mathcal{N}} [h_2(x) \cdot p(x)^k]}_{:= b_k}$$

$$= b_I \cdot \sqrt{2}\alpha_2 + \sum_{k=I}^{C_q} b_k \cdot \varsigma_{k-1}(\alpha), \quad (\text{C.8})$$

694 where $b_I \neq 0$ according to the definition of I . Now we take a closer look at the polynomials
695 $\{\varsigma_k(\alpha_1, \dots, \alpha_{C_q})\}_{k=I-1}^{C_q-1}$. We claim that: (i) they are different polynomials and are linearly independent,
696 and (ii) especially, they are all nonzero. To see these, recall that

$$\varsigma_k(\alpha_1, \dots, \alpha_{C_q}) = \sum_{j=k+1}^{C_q} \sqrt{j} \cdot \alpha_j \cdot r_{j-1,k} \cdot \mathbb{E}_{x \sim \mathcal{N}} \left[x \cdot h_{j-k-1}(x) \cdot (\sigma'(x))^k \right],$$

697 where $\sigma'(x) = \sum_{j=1}^{C_q} \sqrt{j} \cdot \alpha_j \cdot h_{j-1}(x)$. We can calculate that for $k = 0$, $\varsigma_0(\alpha_1, \dots, \alpha_{C_q}) = \sqrt{2}\alpha_2$
698 which is a non-zero polynomial. For $k = 1$, we have that

$$\begin{aligned} \varsigma_1(\alpha_1, \dots, \alpha_{C_q}) &= \sum_{j=2}^{C_q} \sqrt{j} \cdot \alpha_j \cdot r_{j-1,1} \cdot \mathbb{E}_{x \sim \mathcal{N}} [x \cdot h_{j-2}(x) \cdot \sigma'(x)] \\ &= 2r_{1,1} \cdot \alpha_2^2 + \sum_{j=3}^{C_q} \sqrt{j} \cdot \alpha_j \cdot r_{j-1,1} \cdot \mathbb{E}_{x \sim \mathcal{N}} [x \cdot h_{j-2}(x) \cdot \sigma'(x)] \end{aligned}$$

699 which is non-zero (since in the summation from $j = 3$ to C_q there would be no term in the form of
700 c_2^2) and is linearly independent of ς_0 because each terms in the summation here has degree exactly
701 2. Now consider for $k \geq 2$,

$$\begin{aligned} \varsigma_k(\alpha_1, \dots, \alpha_{C_q}) &= \sqrt{2(k+1)k} \cdot r_{k,k} \cdot \alpha_1^{k-1} \alpha_2 \alpha_{k+1} \\ &\quad + \sum_{j=k+2}^{C_q} \sqrt{j} \alpha_j r_{j-1,k} \cdot \mathbb{E}_{x \sim \mathcal{N}} \left[x \cdot h_{j-k-1}(x) \cdot (\sigma'(x))^k \right]. \end{aligned}$$

702 Again, this polynomial is non-zero (since in the summation from $j = k + 2$ to C_q there would be no
703 term in the form of $\alpha_1^{k-1} \alpha_2 \alpha_{k+1}$) and is linearly independent of $\varsigma_0, \dots, \varsigma_{k-1}$ due to the fact that the
704 highest degree of these polynomials is no larger than k , and the order for each term in ς_k is exactly
705 $k + 1$. Thus we have proved the two claims by induction. Now recall that we are aiming at proving the
706 RHS of (C.8) is non-zero. By our two claims just proved, the RHS of (C.8) is a linear combination of
707 $C_q - I + 2$ linearly independent and non-zero polynomials where at least one of the combination
708 coefficient is non-zero (which is b_I). Thus we obtain that the RHS of (C.5) is a non-zero polynomial
709 of $(\alpha_1, \dots, \alpha_{C_q})$ and its zeros form a zero-measure set. This proves that with probability 1 over the
710 randomness of $(\alpha_1, \dots, \alpha_{C_q})$, the high-pass condition holds.

711 **Case 2: $s^* = 1$.** For this case of $s^* = 1$, the proof is almost the same as that for $s^* = 2$, where we
712 additionally utilize the fact that polynomial link function with generative exponent $s^* = 1$ can not be
713 an even polynomial (Example 2.2) and thus there always exists some $I \leq C_q \in \mathbb{N}_+$ such that the
714 information exponent $q^*(p^I) = 1$ (see Proposition 5 of Lee et al. (2024)). With this fact, repeating
715 the above argument can give the desired high-pass property.

716 **Verifying the event \mathcal{E} .** Now we verify that the desired event

$$\mathcal{E} = \{ |\text{err}_{m,l,i}^{(t)}| \leq d^{-10s^*}, \forall (m, l, i, t) \in [M] \times [L] \times [n] \times [T] \}$$

717 holds with probability at least $1 - O(d^{-c_0})$ for some constant $c_0 > 0$ that we specify later. With (C.3)
718 and (C.4), it suffices to look at each of the error terms $\text{err}_{m,l,i,1}^{(t)}$, $\text{err}_{m,l,i,2}^{(t)}$, and $\text{err}_{m,l,i,3}^{(t)}$ respectively.
719 For $\text{err}_{m,l,i,1}^{(t)}$,

$$\left| \text{err}_{m,l,i,1}^{(t)} \right| \leq \sum_{m'=1}^M |a_{m'}| \cdot \left| \sigma(\langle w_{m',l}^{(t)}, z_i \rangle) \right| \cdot \left| \sigma'(\langle w_{m,l}^{(t)}, z_i \rangle) \right| \quad (\text{C.9})$$

720 Note that $\{\langle w_{m,l}^{(t)}, z_i \rangle\}_{m \in [M]}$ are standard Gaussians since $\{w_{m,l}^{(t)}\}_{m \in [M]} \subset \mathbb{S}^{d-1}$. Therefore, with
721 probability at least $1 - d^{-c_0}$ for some constant $c_0 > 0$, we have that $|\langle w_{m,l}, z_i \rangle| = \tilde{O}(1)$. Meanwhile,

722 since that σ and σ' are both polynomials with constant order and bounded coefficients, and that
 723 $M = O(d)$, then by taking $a_m = d^{-11s^*}$, we conclude from (C.9) that

$$|\text{err}_{m,l,i,1}^{(t)}| \leq \tilde{O}(d^{-10s^*}). \quad (\text{C.10})$$

724 For the second error term $\text{err}_{m,l,i,2}^{(t)}$, similarly we have that

$$\left| \text{err}_{m,l,i,2}^{(t)} \right| \leq \sum_{m' \in [M]} |a_{m'}| \cdot \left| \sigma(\langle (w_{m',l}^{(t)})^+, z_i \rangle) \right| \cdot \left| \sigma'(\langle (w_{m,l}^{(t)})^+, z_i \rangle) \right|. \quad (\text{C.11})$$

725 Note that the one-step updated weights satisfy that

$$\left| \langle (w_{m,l}^{(t)})^+, z_i \rangle \right| = \left| \langle w_{m,l}^{(t)}, z_i \rangle + y_i \cdot \sigma'(\langle w_{m,l}^{(t)}, z_i \rangle) + \text{err}_{m,l,i,1}^{(t)} \right| = \tilde{O}(1), \quad (\text{C.12})$$

726 with probability at least $1 - d^{-c_0}$ since $y_i = p(\langle \theta^*, z_i \rangle)$ and p is also a polynomial of constant degree
 727 and coefficients. Therefore, with the choice of a_m 's, by (C.11), we conclude that

$$\left| \text{err}_{m,l,i,2}^{(t)} \right| \leq \tilde{O}(d^{-10s^*}). \quad (\text{C.13})$$

728 Finally, regarding $\text{err}_{m,l,i,3}^{(t)}$, note that with the same argument as (C.12), we know that with probability
 729 at least $1 - d^{-c_0}$, both $\langle w_{m,l}^{(t)}, z_i \rangle + y_i \cdot \sigma'(\langle w_{m,l}^{(t)}, z_i \rangle) + \text{err}_{m,l,i,1}^{(t)}$ and $\langle w_{m,l}^{(t)}, z_i \rangle + y_i \cdot \sigma'(\langle w_{m,l}^{(t)}, z_i \rangle)$
 730 are $\tilde{O}(1)$. Since σ' is a polynomial, it is $\tilde{O}(1)$ -Lipschitz continuous for inputs that are $\tilde{O}(1)$. Therefore,
 731 combined with (C.10) that we have proved, we can obtain that

$$\left| \text{err}_{m,l,i,3}^{(t)} \right| \leq |y_i| \cdot \tilde{O}(|\text{err}_{m,l,i,1}^{(t)}|) = \tilde{O}(d^{-10s^*}). \quad (\text{C.14})$$

732 Finally, combining (C.10), (C.13), and (C.14), we obtain that for given (m, l, i, t) , with probability at
 733 least $1 - O(d^{-c_0})$, it holds that

$$|\text{err}_{m,l,i}^{(t)}| = \tilde{O}(d^{-10s^*}).$$

734 Finally, taking c_0 as a constant that is larger than 2 and applying a union bound argument, we can
 735 obtain that

$$\Pr(\mathcal{E}) \geq 1 - MLnT \cdot \tilde{O}(d^{-c_0}) \geq 1 - \tilde{\Theta}(d^2) \cdot \tilde{O}(d^{-c_0}) \geq 1 - \tilde{O}(d^{-c'_0})$$

736 for some other constant $c'_0 > 0$. Here we have applied our choice of (M, L, n, T) in our algorithm
 737 (see Algorithm 1, Algorithm 2). Thus we verify the property of the event \mathcal{E} , proving Corollary 4.5. \square

738 **C.2.2 Modified loss for general $s^* \geq 1$**

739 Here we give a specific choice of the activation function σ and the loss function ℓ . We mainly focus on
 740 the situation where \mathbb{Q}_y has a continuous cumulative distribution function $F_{\mathbb{Q}_y}$ with bounded density
 741 $f_{\mathbb{Q}_y}$. For the situation where \mathbb{Q}_y is a discrete distribution (e.g., classification task), we discuss them in
 742 the end of this section. For the activation function σ , we let $\sigma(x) := (1/\sqrt{s^*}) \cdot h_{s^*}(x)$. Since then,

$$\hat{\psi}_s(y) = \mathbb{E}_{\mathbb{Q}}[\psi(x, y) \cdot h_s(x) | y] = \mathbb{E}_{\mathbb{Q}}[\sigma'(x) \cdot h_s(x)] \cdot \ell'(y) = 0, \quad \forall s < s^* - 1. \quad (\text{C.15})$$

743 Regarding the choice of the loss function ℓ , we remark that if one chooses a fixed loss function, there
 744 always exist instances such that the second assumption in the high-pass condition fails. To address
 745 this issue, we propose to construct a random loss function ℓ . To rule out pathological examples of the
 746 underlying distribution \mathbb{P} , we make the following assumption on the coefficient function ζ_{s^*} .

747 **Assumption C.2.** We assume that the expansion of $\tilde{\zeta}_{s^*} := \zeta_{s^*} \circ F_{\mathbb{Q}_y}^{-1} : [0, 1] \mapsto \mathbb{R}$ on the Fourier
 748 basis $\{\varphi_i(x)\}_{i \geq 0}$ of $[0, 1]$ has a non-zero coefficient of order at most $D = O(1)$.

749 We then choose the loss function ℓ as the following,

$$\ell'(y) = \sum_{i=0}^D \alpha_i \cdot \varphi_i \circ F_{\mathbb{Q}_y}(y), \quad \alpha_i \sim \text{Unif}([0, 1]), \quad \forall 0 \leq i \leq D.$$

750 Notice that $F_{\mathbb{Q}_y}$ can be estimated from data using a one dimensional density estimator. Thus here we
 751 directly assume the accessibility of the function $F_{\mathbb{Q}_y}$. This further gives that

$$\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] = \mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \ell'(y)] = \sum_{i=0}^D \alpha_i \cdot \mathbb{E}_{\text{Unif}([0,1])} [\tilde{\zeta}_{s^*}(\tilde{y}) \cdot \varphi_i(\tilde{y})], \quad (\text{C.16})$$

752 which is a non-zero polynomial of the coefficients $\{\alpha_i\}_{i \leq D}$ due to Assumption C.2.

753 *Proof of Corollary 4.7.* To prove Corollary 4.7, it suffices to show that (i) Assumption 4.1 holds, and
 754 (ii) the event \mathcal{E} holds with the desired high probability. In the following, we first verify Assumption 4.1,
 755 and then check the event \mathcal{E} .

756 **Verifying Assumption 4.1.** First, since σ' is a polynomial and ℓ' is bounded (the Fourier basis is
 757 bounded and $D = O(1)$), we know that both Assumption 4.1(a) and Assumption 4.1(c) are satisfied.
 758 Then, by the discussions before the proof, we know from (C.15) that the first condition in the high-pass
 759 assumption (Assumption 4.1(b)) is satisfied. Furthermore, according to (C.16) and Assumption C.2,
 760 $\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)]$ is a non-zero polynomial of the coefficients $\{\alpha_i\}_{i \leq D}$ and thus its zeros form
 761 a measure-zero set. This means that with probability 1 over the randomness of $(\alpha_1, \dots, \alpha_D)$, the
 762 second condition in the high-pass assumption is also satisfied. This verifies Assumption 4.1.

763 **Verifying the event \mathcal{E} .** Recall our definition in Example 4.6, the error term is defined as

$$\text{err}_{m,l,i}^{(t)} = \left(\ell'(y_i) - \ell' \left(y_i - f(z_i; \{w_{m,l}^{(t)}\}_{m \in [M]}) \right) \right) \cdot \sigma'(\langle w_{m,l}^{(t)}, z_i \rangle).$$

764 First, since $w_{m,l}^{(t)} \in \mathbb{S}^{d-1}$, $\langle w_{m,l}^{(t)}, z_i \rangle$ is a standard Gaussian and therefore $|\langle w_{m,l}^{(t)}, z_i \rangle| = \tilde{O}(1)$ with
 765 probability at least $1 - d^{-c_0}$ for some constant $c_0 > 0$. Since σ' is a polynomial of constant degree,
 766 we then obtain that $|\sigma'(\langle w_{m,l}^{(t)}, z_i \rangle)| = \tilde{O}(1)$ with probability at least $1 - d^{-c_0}$. Second, consider that
 767 the second derivative of the loss function $\ell''(y)$ is given by

$$\ell''(y) = \sum_{i=0}^D \alpha_i \cdot \varphi_i'(F_{\mathbb{Q}_y}(y)) \cdot f_{\mathbb{Q}_y}(y),$$

768 which satisfies $|\ell''(y)| = O(1)$ since the derivative of the Fourier basis is still bounded and that the
 769 density of \mathbb{Q}_y is assumed to be bounded. Therefore, we have

$$\begin{aligned} \ell'(y_i) - \ell' \left(y_i - f(z_i; \{w_{m,l}^{(t)}\}_{m \in [M]}) \right) &= O \left(\left| f(z_i; \{w_{m,l}^{(t)}\}_{m \in [M]}) \right| \right) \\ &= O \left(\sum_{m=1}^M |a_m| \cdot \left| \sigma(\langle w_{m,l}^{(t)}, z_i \rangle) \right| \right). \end{aligned}$$

770 Since σ is a polynomial of constant degree and $\langle w_{m,l}^{(t)}, z_i \rangle$ are all standard Gaussians, we have that
 771 $|\sigma(\langle w_{m,l}^{(t)}, z_i \rangle)| = \tilde{O}(1)$ with probability at least $1 - O(d^{-c_0})$. Now given that $M = O(d)$ and taking
 772 $a_m = d^{-11s^*}$, we have that

$$\ell'(y_i) - \ell' \left(y_i - f(z_i; \{w_{m,l}^{(t)}\}_{m \in [M]}) \right) = \tilde{O}(d^{-10s^*}),$$

773 with probability at least $1 - O(d^{-c_0})$. Therefore, for any given (m, l, i, t) , with probability at least
 774 $1 - \tilde{O}(d^{-c_0})$, it holds that

$$|\text{err}_{m,l,i}^{(t)}| = \tilde{O}(d^{-10s^*}) \cdot \tilde{O}(1) = \tilde{O}(d^{-10s^*}).$$

775 Finally, as in the proof of Corollary 4.5, we take c_0 as a constant that is larger than $s^* + 1$, and apply
 776 a union bound argument, by which we can obtain that

$$\Pr(\mathcal{E}) \geq 1 - MLnT \cdot \tilde{O}(d^{-c_0}) \geq 1 - \tilde{\Theta}(d^{s^*+1/2}) \cdot \tilde{O}(d^{-c_0}) \geq 1 - \tilde{O}(d^{-c'_0})$$

777 for some other constant $c'_0 > 0$. Thus we verify the property of the event \mathcal{E} , proving Corollary 4.7. \square

778 **Remark C.3** (Discrete labels). For the case of discrete label y that supports on a finite set \mathcal{Y} (e.g.,
779 classification tasks), the construction of ψ is somehow more direct. In this case, we can still consider
780 an oracle function in the form of $\psi(y, x) = \sigma'(x) \cdot \varphi(y)$ for some function $\varphi(y)$. The activation
781 function $\sigma(x) = (1/\sqrt{s^*}) \cdot h_{s^*}(x)$ and the function $\varphi(y)$ is a random function given by

$$\varphi(y) \sim \text{Unif}([0, 1]), \quad \forall y \in \mathcal{Y}. \quad (\text{C.17})$$

782 On the one hand, we can directly conclude as in (C.15) that $\hat{\psi}_s(y) = 0$ for all $s < s^* - 1$. On the
783 other hand, we have that

$$\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] = \mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \varphi(y)] = \sum_{y \in \mathcal{Y}} \mathbb{Q}_y(y) \cdot \zeta_{s^*}(y) \cdot \varphi(y).$$

784 By the definition of generative exponent (Definition 2.1), $\mathbb{E}_{\mathbb{Q}_y}[\zeta_{s^*}(y)^2] > 0$ and thus at least one of
785 $\{\mathbb{Q}_y(y) \cdot \zeta_{s^*}(y)\}_{y \in \mathcal{Y}}$ is non-zero. Thus under (C.17), $\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)]$ is non-zero with proba-
786 bility 1 over the randomness in φ . Thus we have verified Assumption 4.1(b). Assumption 4.1(a) and
787 Assumption 4.1(c) can be verified in the same way as for the continuous case, and thus Assumption 4.1
788 is checked. Finally, we remark that in the discrete case we do not attempt to reduce the oracle function
789 from certain loss derivative and thus we simply set the error terms err as zero. Thus all the conditions
790 in Theorem 4.2 hold and Corollary 4.7 is proved.

791 D Proof Sketch of the Main Theorem for Uniform Prior

792 For simplicity, denote $\rho_m^{(t)} := \langle \theta_m^{(t)}, \theta^* \rangle$, the alignment between the weights of neuron m and the
793 signal θ^* at time t . Recall from Line 8 in Algorithm 1 that the update for neuron m at time step t is

$$\theta_m^{(t+1)} = \frac{\theta_m^{(t)} + \eta \bar{g}_m^{(t)}}{\|\theta_m^{(t)} + \eta \bar{g}_m^{(t)}\|_2}. \quad (\text{D.1})$$

794 This implies that the alignment of the next iteration, $\rho_m^{(t+1)}$, is a convex combination of the previous
795 alignment $\rho_m^{(t)}$ and the alignment of the update step $\langle \bar{g}_m^{(t)}, \theta^* \rangle = \langle g_m^{(t)}, \theta^* \rangle / \|g_m^{(t)}\|_2$. Therefore, to
796 show that the alignment improves after one iteration, we need to first analyze the scale of $\langle g_m^{(t)}, \theta^* \rangle$
797 and $\|g_m^{(t)}\|_2$; then we will be able to characterize the improvement of $\rho_m^{(t)}$ across iterations.

798 **Alignment of the update step $\langle \bar{g}_m^{(t)}, \theta^* \rangle$.** To this end, we calculate the first moment and second
799 moment of $g_m^{(t)}$ over the randomness of the data $\{(z_i^{(t)}, y_i^{(t)})\}_{i=1}^n$, and combining these leads to the
800 concentration of $\langle g_m^{(t)}, \theta^* \rangle / \|g_m^{(t)}\|_2$. More specifically, for the first moment of $g_m^{(t)}$, we have

$$\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g_m^{(t)}], \theta^* \rangle \approx \rho_m^{(t)} \gamma \cdot (|\rho_m^{(t)}| \gamma + d^{-1/2})^{s^*-2},$$

801 while the magnitude of $\mathbb{E}_{\mathbb{P}_{\theta^*}}[g_m^{(t)}]$ in any other direction orthogonal to θ^* is of strictly higher order. For
802 the second moment of $g_m^{(t)}$, setting $\gamma = \tilde{\Theta}(d^{-1/4})$, it can be shown that for any direction $v \in \mathbb{S}^{d-1}$,
803 $\mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_m^{(t)}, v \rangle^2] = \tilde{O}(d^{-(s^*-1)/2})$. Now choosing $n = \tilde{\Omega}(d^{s^*/2})$, it follows from Bernstein-type
804 concentration inequality that the fluctuation of $\langle g_m^{(t)}, v \rangle$ is of the same order $\tilde{O}(d^{-s^*/2+1/4})$ for any
805 direction $v \in \mathbb{S}^{d-1}$. Therefore, with high probability,

$$\begin{aligned} \langle g_m^{(t)}, \theta^* \rangle &\geq \rho_m^{(t)} \gamma \cdot (|\rho_m^{(t)}| \gamma + d^{-1/2})^{s^*-2} - \tilde{O}(d^{-s^*/2+1/4}), \\ \|g_m^{(t)}\|_2 &\leq \rho_m^{(t)} \gamma \cdot (|\rho_m^{(t)}| \gamma + d^{-1/2})^{s^*-2} + \tilde{O}(d^{-s^*/2+3/4}). \end{aligned} \quad (\text{D.2})$$

806 **Phase 1: from $d^{-1/2}$ to weak alignment.** Due to random initialization, it holds with high proba-
807 bility that $\rho_m^{(0)} = O(d^{-1/2})$ for $\Omega(M)$ many neurons. Therefore, it suffices to consider a neuron
808 with $\rho_m^{(t)} = \Omega(d^{-1/2})$. When $\Omega(d^{-1/2}) \leq \rho_m^{(t)} \leq O(1)$, by choosing $\gamma = \tilde{\Theta}(d^{-1/4})$, we can ensure
809 that the first term in the lower bound for $\langle g_m^{(t)}, \theta^* \rangle$ in (D.2) dominates the $\tilde{O}(d^{-s^*/2+1/4})$ fluctuation.
810 Based on this, we can leverage (D.2) to further show that $\langle \bar{g}_m^{(t)}, \theta^* \rangle = \langle g_m^{(t)}, \theta^* \rangle / \|g_m^{(t)}\|_2 \geq (1+c)\rho_m^{(t)}$
811 for a constant $c > 0$. Consequently, it follows from (D.1) that

$$\rho_m^{(t+1)} \geq (1+c)\rho_m^{(t)} \quad \text{for some constant } c > 0.$$

812 Therefore, it takes $O(\log d)$ many steps for $\rho_m^{(t)}$ to increase from $d^{-1/2}$ to $O(1)$. During this period,
813 the dynamics will go through two phases separated by a critical alignment level ρ^* such that

$$\rho_m^{(t)} \gamma \cdot (|\rho_m^{(t)}| \gamma + d^{-1/2})^{s^* - 2} \approx \tilde{O}(d^{-s^*/2 + 3/4}),$$

814 which gives that $\rho^* = \tilde{\Theta}(d^{-1/4})$. After the alignment $\rho_m^{(t)}$ reaches ρ^* , there is a short period where
815 $\rho_m^{(t)}$ grows rapidly as $\rho_m^{(t+1)}/\rho^* \geq (\rho_m^{(t)}/\rho^*)^{s^* - 1}$, until the alignment reaches $d^{-1/4 + 1/4(s^* - 1)}$. The
816 length of this period is very short compared to the other periods on the road to weak alignment.

817 **Phase 2: from weak to strong alignment.** Finally, after $\rho_m^{(t)}$ grows to a constant scale, we need to
818 track the value of $1 - \rho_m^{(t)}$. Again using (D.2), we can show that $1 - \rho_m^{(t+1)} \leq (1 + c)(1 - \rho_m^{(t)})$ for
819 some constant $c > 0$. Hence, it takes another $O(\log d)$ steps to eventually achieve strong alignment.

820 E Proof of the Main Theorem for the Uniform Prior

821 Now we present the proof for Theorem 4.2. We introduce a shorthand $\rho = \langle \theta, \theta^* \rangle$ for the alignment
822 between θ and θ^* . This shorthand inherits the subscript and superscript of θ as well, i.e., $\rho_m^{(t)} =$
823 $\langle \theta_m^{(t)}, \theta^* \rangle$.

824 Recall from Algorithm 1 that at the t -th step, given the normalized gradient step $g_m^{(t)} = g_m^{(t)} / \|g_m^{(t)}\|_2$,
825 the updated weight parameter is given by

$$\theta_m^{(t+1)} = \frac{\theta_m^{(t)} + \eta g_m^{(t)} / \|g_m^{(t)}\|_2}{\|\theta_m^{(t)} + \eta g_m^{(t)} / \|g_m^{(t)}\|_2\|_2}.$$

826 Note that the alignment $\langle \theta_m^{(t+1)}, \theta^* \rangle$ depends on the alignment of the previous iterate $\langle \theta_m^{(t)}, \theta^* \rangle$ and
827 the alignment of the current update step $\langle g_m^{(t)} / \|g_m^{(t)}\|_2, \theta^* \rangle$, so we first need to analyze the latter.

828 Here, we stop to introduce an immediate result that is crucial in characterizing the alignment.

829 **Almost orthogonality of smoothing noise.** Recall that the perturbed weights are $w_{m,l}^{(t)} =$
830 $(\gamma \theta_m^{(t)} + \xi_{m,l}) / \|\gamma \theta_m^{(t)} + \xi_{m,l}\|_2$, where $\xi_{m,l} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1})$. Due to the high dimensionality, $\xi_{m,l}$ is
831 almost orthogonal to any designated direction with high probability. In the context, we are primarily
832 interested in the alignment with θ^* . Correspondingly, we decompose $g_m^{(t)}$ with respect to the following
833 orthonormal basis

$$\{v_{m,1}^{(t)} = \theta^*, v_{m,2}^{(t)} = (1 - \rho^2)^{-1/2} \cdot (\theta_m^{(t)} - \langle \theta_m^{(t)}, \theta^* \rangle \cdot \theta^*), v_{m,3}^{(t)}, \dots, v_{m,d}^{(t)}\} \quad (\text{E.1})$$

834 We justify this property by defining the following nice event for $\epsilon > 0$:

$$\mathcal{E}_m^{(t)}(\epsilon) = \left\{ \max_{1 \leq i \leq d} |\langle \xi_{m,l}, v_{m,i}^{(t)} \rangle| < \epsilon, \forall l \in [L] \right\}.$$

835 This event helps in characterizing the alignment between the expected gradient $\mathbb{E}_{\mathbb{P}_{\theta^*}}[g_m^{(t)}]$ and the
836 signal θ^* . Additionally, we define another event for $\tilde{\epsilon} > 0$:

$$\tilde{\mathcal{E}}_m^{(t)}(\tilde{\epsilon}) = \{|\langle \xi_{m,l}, \xi_{m,l'} \rangle| < \tilde{\epsilon}, \forall l, l' \in [L] \text{ s.t. } l \neq l'\}.$$

837 This event controls the correlation between the noise vectors, which helps in controlling the fluctuation
838 of the gradient. The following lemma provides some direct benefits of these events.

839 **Lemma E.1** (Polarized weight on the nice event). *For the orthonormal directions*
840 $\{\theta^*, v_{m,2}^{(t)}, \dots, v_{m,d}^{(t)}\}$ *defined in (E.1), suppose that the corresponding nice event* $\mathcal{E}_m^{(t)}(\epsilon) \cap \tilde{\mathcal{E}}_m^{(t)}(\tilde{\epsilon})$
841 *holds and the polarization level* $\gamma \in (0, 1/2)$. *Then we have for any* $l \in [L]$ *that*

$$|\langle w_{m,l}^{(t)}, \theta^* \rangle| \leq 2(\gamma |\langle \theta_m^{(t)}, \theta^* \rangle| + \epsilon), \quad |\langle w_{m,l}^{(t)}, v_{m,1}^{(t)} \rangle| \leq 2\left(\gamma \sqrt{1 - \langle \theta_m^{(t)}, \theta^* \rangle^2} + \epsilon\right),$$

$$|\langle w_{m,l}^{(t)}, v_{m,i}^{(t)} \rangle| \leq 2 \cdot \epsilon, \quad 2 \leq i \leq d.$$

842 Additionally, for any $l \neq l'$, we have that

$$\langle w_{m,l}^{(t)}, w_{m,l'}^{(t)} \rangle \leq 4(\gamma^2 + 2\gamma\epsilon + \tilde{\epsilon}).$$

843 *Proof of Lemma E.1.* See Appendix E.3. □

844 **Characterizing** $\langle g_m^{(t)}, \theta^* \rangle$. Note that $g_m^{(t)} = n^{-1} \sum_i g_{m,i}^{(t)}$, where

$$g_{m,i}^{(t)} = \frac{1}{L} \sum_{l=1}^L (\psi(y_i^{(t)}, \langle w_{m,l}^{(t)}, z_i^{(t)} \rangle) \cdot z_i^{(t)} - \hat{\psi}_1(y_i^{(t)}) \cdot w_{m,l}).$$

845 We characterize the alignment $\langle g_m^{(t)}, \theta^* \rangle / \|g_m^{(t)}\|_2$ in Appendix E.1 with two steps, stated in two key
846 propositions as follows:

- 847 1. In Proposition E.3, we analyze the magnitude of $\mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_m^{(t)}, \theta^* \rangle]$ and $\|\mathbb{E}_{\mathbb{P}_{\theta^*}}[P_{\theta^*}^\perp \langle g_m^{(t)}, \theta^* \rangle]\|$.
848 2. In Proposition E.4, we control the fluctuation of $\langle g_m^{(t)}, \theta^* \rangle$ around its expectation using the
849 polynomial-tail like property in Assumption 4.1(c).

850 Both propositions are established under the nice event $\mathcal{E}_m^{(t)}(\epsilon) \cap \tilde{\mathcal{E}}_m^{(t)}(\tilde{\epsilon})$. The proof of these propo-
851 sitions is deferred to Appendix E.3. Finally, with these two propositions, we prove Theorem 4.2 in
852 Appendix E.2.

853 E.1 Properties of the Gradient Step

854 In this part, we characterize the alignment of normalized update $g_m^{(t)} / \|g_m^{(t)}\|$ with the signal θ^* , given
855 $\theta_m^{(t)}$. Since we are focusing on the one step behavior for a fixed neuron $m \in [M]$, we omit the neuron
856 index m and time index t in the sequel. To facilitate the presentation, we propose the following
857 simplified setup that extract all the essential elements to describe the one-step behavior.

858 **Definition E.2.** Fix θ and θ^* and let $\rho = \langle \theta, \theta^* \rangle$. Suppose the data points $(z_1, y_1), \dots, (z_n, y_n)$ are
859 i.i.d. generated from \mathbb{P}_{θ^*} . Define $w_l = (\gamma\theta + \xi_l) / \|\gamma\theta + \xi_l\|_2$ for $l = 1, \dots, L$, where ξ_1, \dots, ξ_L $\stackrel{\text{i.i.d.}}{\sim}$
860 $\text{Unif}(\mathbb{S}^{d-1})$ are independent of $\{(z_i, y_i)\}_{i=1}^n$. Given the oracle $\psi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, we define

$$g = \frac{1}{nL} \sum_{i=1}^n \sum_{l=1}^L (\psi(y_i, \langle w_l, z_i \rangle) \cdot z_i - \hat{\psi}_1(y_i) \cdot w_l).$$

861 To describe the associated good event, we fix an orthonormal basis:

$$v_1 = \theta^*, v_2 = \frac{\theta - \rho\theta^*}{\sqrt{1 - \rho^2}}, v_3, \dots, v_d,$$

862 and define

$$\mathcal{E}(\epsilon) = \left\{ |\langle \xi_l, \theta^* \rangle| < \epsilon, \quad \max_{2 \leq i \leq d} |\langle \xi_l, v_i \rangle| < \epsilon, \quad \forall l \in [L] \right\};$$

$$\tilde{\mathcal{E}}(\tilde{\epsilon}) = \{ |\langle \xi_l, \xi_{l'} \rangle| < \tilde{\epsilon}, \quad \forall l, l' \in [L] \text{ s.t. } l \neq l' \}.$$

863 As mentioned in Appendix D, we can reduce this problem to first characterizing $\mathbb{E}_{\mathbb{P}_{\theta^*}}[g]$, and then
864 control the fluctuation of g around its expectation. To this end, we first introduce a lemma that
865 characterizes the first moment of the gradient step. This lemma is valid for both the non-sparse and
866 sparse setting and is helpful in understanding the structure of the expected gradient.

867 **Lemma H.2** (Decomposition of the first moment). Suppose that we are working with the setting in
868 Definition E.2, where the oracle function ψ follows Assumption 4.1. Then it holds that

$$\mathbb{E}_{\mathbb{P}_{\theta^*}}[g] = \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s-1}(y)] \cdot \frac{\sqrt{s}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^{s-1} \cdot \theta^*$$

$$+ \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)] \cdot \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^s \cdot w_l.$$

869 *Proof of Lemma H.2.* See Appendix H. □

870 One can easily see that the first summation term corresponds to the signal, while the second corre-
871 sponds to the resilience of current weight. The structure of w_l guarantees $\langle w_l, \theta^* \rangle$ is small, therefore
872 only the leading term in the geometric series above is dominant. Also we note that the leading term
873 from the signal is larger than the leading term from the resilience, which indicates that the expected
874 gradient is highly aligned with the true signal. This is justified in the next proposition.

875 Before stating it, we fix M to be a sufficiently large constant that does not scale with d and $T =$
876 $O(\log d)$. The involvement of M, T here is merely for the union bound argument in Appendix E.2.

877 **Proposition E.3** (Alignment of expected gradient). *Suppose that we are working with the setting in*
878 *Definition E.2, where the oracle function ψ follows Assumption 4.1. Additionally, we set $\gamma = o(1)$,*
879 *$L = \Omega((\epsilon \vee \gamma)^{s^* - 1} \cdot d^{s^*/2} \vee (d \log d))$ and $\epsilon = o(1)$ is chosen such that $\Pr(\mathcal{E}(\epsilon)) = 1 - O(d^{-s^*/2})$.*
880 *Then there exists a $\{\xi_l\}_{l \in [L]}$ -measurable event \mathcal{E}_1 with $\Pr(\mathcal{E}_1) \geq 1 - d^{-c}(MT)^{-1}$, such that on the*
881 *event $\mathcal{E}_1 \cap \mathcal{E}(\epsilon)$, it holds that*

$$\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g], \theta^* \rangle \simeq \begin{cases} \gamma \rho \cdot (\gamma |\rho| + d^{-1/2})^{s^* - 2} & \text{if } s^* \text{ is even;} \\ (\gamma |\rho| + d^{-1/2})^{s^* - 1} & \text{if } s^* \text{ is odd,} \end{cases}$$

882 and that

$$\left| \|\mathbb{E}_{\mathbb{P}_{\theta^*}}[g]\|_2 - |\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g], \theta^* \rangle| \lesssim (\gamma |\rho| + d^{-1/2})^{s^*},$$

883 as long as $\gamma |\rho| = \omega(d^{-1})$.

884 *Proof of Proposition E.3.* See Appendix E.3. □

885 From this proposition, it is already clear that the expected gradient is highly aligned with the signal
886 θ^* in the sense that $\|P_{\theta^*}^\perp \mathbb{E}_{\mathbb{P}_{\theta^*}}[g]\|_2 < |\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g], \theta^* \rangle|$ whenever $\gamma |\rho| = \omega(d^{-1})$. Later we will see
887 that this is indeed the case during the trajectory of Algorithm 1.

888 **Proposition E.4** (Fluctuation of mini-batch gradient). *Under the simplified setting introduced in*
889 *Definition E.2 where $\psi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ follows Assumption 4.1. Suppose that we choose ϵ and $\tilde{\epsilon}$ such*
890 *that*

$$\epsilon^2 \leq \tilde{\epsilon} \ll 1; \quad 2\gamma\epsilon \leq \tilde{\epsilon}.$$

891 Also, suppose that sample size

$$n = \Omega\left(\left((\gamma^2 + \tilde{\epsilon})^{s^* - 1} + L^{-1}\right)^{-1} \cdot \log(d)^{2C_p + 2}\right)$$

892 where C_p is defined in Assumption 4.1(c). Then there exists a $\{(y_i, z_i)\}_{i \in [n]}$ -measurable event \mathcal{E}_2
893 with $\Pr(\mathcal{E}_2^c) \leq d^{-c} \cdot (MT)^{-1}$. And it holds on $\mathcal{E}_2 \cap \mathcal{E}(\epsilon) \cap \tilde{\mathcal{E}}(\tilde{\epsilon})$ that

$$|\langle g, \theta^* \rangle - \langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g], \theta^* \rangle| \lesssim \sqrt{\frac{\left((\gamma^2 + \tilde{\epsilon})^{s^* - 1} + L^{-1}\right) \cdot \log(d)}{n}},$$

894 and that

$$\|g - \mathbb{E}_{\mathbb{P}_{\theta^*}}[g]\|_2 \lesssim \sqrt{\frac{\left((\gamma^2 + \tilde{\epsilon})^{s^* - 1} + L^{-1}\right) \cdot d \log(d)}{n}}.$$

895 *Proof of Proposition E.4.* See Appendix E.3. □

896 E.2 Proof of the Main Theorem for Uniform Prior

897 Now we are ready to present the proof of Theorem 4.2.

898 *Proof of Theorem 4.2.* We first establish the good events required for the proof, characterize the
899 properties of the update step on these events, and then put things together to establish the alignment
900 of the model weights with the signal.

901 **Preparations.** To start with, we clarify the event we will work with by verifying that our configu-
 902 ration is compatible with the conditions in Proposition E.3 and Proposition E.4. Recall that we set
 903 $L = \Omega(d^{(s^*-1)/2} \vee (d \log d)^{1/2})$ and is at most polynomial in d , the scale of L clearly satisfy that
 904 $L = \Omega((d^{1/2} \epsilon)^{s^*} \vee (d \log d))$. During our algorithm, $\gamma = (d^{-1} \cdot \log d)^{1/4} = o(1)$ is fixed. Choosing
 905 $\epsilon = d^{-1/2} \log d$, we have by Lemma J.6 that for any t and m , it holds that

$$1 - \Pr(\mathcal{E}_m^{(t)}(\epsilon)) \leq Ld \cdot (\exp(-d/16) + d^{-\log d/4}),$$

906 which decays faster than any constant-degree polynomial in d . Therefore, for sufficiently large d , it
 907 holds that $\Pr(\mathcal{E}_m^{(t)}(\epsilon)^c) = O(d^{-s^*/2})$. So far, we see that all the conditions in Proposition E.3 are
 908 satisfied and we denote the associated event as $\mathcal{E}_{m,1}^{(t)}$.

909 Next, we verify the conditions in Proposition E.4. We choose $\tilde{\epsilon} = \sqrt{4(c + \log_d(MTL^2)) \cdot d^{-1} \log d}$,
 910 then it holds by Lemma J.6 that

$$1 - \Pr(\tilde{\mathcal{E}}_m^{(t)}(\tilde{\epsilon})) \leq L^2 \cdot (\exp(-d/16) + d^{-\tilde{\epsilon}^2 \cdot \log d/4}) \lesssim d^{-c}/MT.$$

911 Additionally, we see that both $\epsilon^2 \leq \tilde{\epsilon} \ll 1$ and $2\gamma\epsilon \leq \tilde{\epsilon}$ is satisfied for sufficiently large d . It is easily
 912 verified that our choice of $L = \Omega(d^{(s^*+1)/2} \vee (d \log d))$ clearly meets the condition that

$$L \gtrsim (\epsilon \vee \gamma)^{s^*-1} \cdot d^{s^*/2} \vee d \log d$$

913 we have that the sample size threshold is now

$$\frac{\log(d)^{2C_p+2}}{((\gamma^2 + \tilde{\epsilon})^{s^*-1} + L^{-1})} \lesssim \log(d)^{2C_p+2} \cdot d^{(s^*-1)/2},$$

914 which is satisfied by our choice $n = \Theta(((d \log d)^{s^*/2} \vee d \log d) \log d)$. Hence, all the conditions in
 915 Proposition E.4 are satisfied and we denote the associated event as $\mathcal{E}_{m,2}^{(t)}$.

916 Recalling that the gradient in Definition E.2 does not include the error term $\text{err}_{m,l,i}^{(t)}$, we additionally
 917 need an event that controls the norm of the inputs $z_{i_2}^{(t)}$, which helps to control $\|\text{err}_{m,l,i}^{(t)} \cdot z_i\|_2$. For
 918 this purpose, we define

$$\mathcal{E}_{m,3}^{(t)} = \left\{ \max_{i \in [n]} \|z_i^{(t)}\|_2 \leq \sqrt{d} \right\}.$$

919 By standard Bernstein's inequality, we have that $\Pr(\mathcal{E}_{m,3}^{(t)}) \leq Ld \cdot \exp\{-d/8\} = O(\exp\{-C'd\})$
 920 for some $C' > 0$. To put things together, we work on the following event:

$$\mathcal{E} = \bigcap_{m=1}^M \bigcap_{t=1}^T (\mathcal{E}_m^{(t)}(\epsilon) \cap \tilde{\mathcal{E}}_m^{(t)}(\tilde{\epsilon}) \cap \mathcal{E}_{m,1}^{(t)} \cap \mathcal{E}_{m,2}^{(t)} \cap \mathcal{E}_{m,3}^{(t)}),$$

921 which is of $\Pr(\mathcal{E}) \geq 1 - O(d^{-c})$ for some $c > 0$ by the union bound argument. Denote

$$\bar{g}_m^{(t)} = \frac{1}{nL} \sum_{i=1}^n \sum_{l=1}^L (\psi(y_i^{(t)}, \langle w_{m,l}^{(t)}, z_i^{(t)} \rangle) \cdot z_i^{(t)} - \hat{\psi}_1(y_i^{(t)}) \cdot w_{m,l}),$$

922 then $\bar{g}_m^{(t)}$ and the mini-batch data $\{(y_i^{(t)}, z_i^{(t)})\}_{i \in [n]}$ match the definition in Definition E.2, which
 923 allows us to apply Proposition E.3 and Proposition E.4. Thanks to the event $\mathcal{E}_{m,3}^{(t)}$, we always have for
 924 any $v \in \mathbb{S}^{d-1}$ that

$$\begin{aligned} |\langle g_m^{(t)}, v \rangle - \langle \bar{g}_m^{(t)}, v \rangle| &\leq \left| \|g_m^{(t)}\|_2 - \|\bar{g}_m^{(t)}\|_2 \right| \\ &\leq d^{1/2} \cdot \max_{l,i} |\text{err}_{m,l,i}^{(t)}| \\ &\leq d^{-9s^*}. \end{aligned} \tag{E.2}$$

925 In the sequel, we restrict our attention to neurons that have $d^{-1/2}/2$ alignment, i.e., the index m such
 926 that $|\langle \theta_m^{(0)}, \theta^* \rangle| \geq d^{-1/2}/2$. From now on, we will drop the neuron index m and the iteration index
 927 (t) in the following analysis for simplicity. The updated weight parameter is denoted as θ' , and the
 928 alignment after the update is denoted as $\rho' = \langle \theta', \theta^* \rangle$. Note that for large $M \gg 1$, the number of
 929 neurons with initial alignment $|\rho| \geq d^{-1/2}/2$ is at least $\Omega(M)$. For our convenience, in the following
 930 we will denote by

$$\kappa := \frac{n}{(d \log d)^{s^*/2} \vee d \log d} \cdot (\log d)^{-1} = \Omega(1).$$

931 Under the preceding configuration, Proposition E.4 and Eq. (E.2) together imply that the fluctuations
 932 of $\langle g, \theta^* \rangle$ can be further bounded by

$$\begin{aligned} |\langle g, \theta^* \rangle - \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}, \theta^* \rangle]| &\lesssim \sqrt{\frac{((\gamma^2 + \tilde{\epsilon})^{s^*-1} + L^{-1}) \cdot \log(d)}{n}} + |\langle \bar{g}, \theta^* \rangle - \langle g, \theta^* \rangle| \\ &\lesssim \sqrt{\frac{(d^{-1} \log d)^{(s^*-1)/2} \cdot \log(d)}{n}} + d^{-9s^*} \\ &= \begin{cases} d^{-(2s^*-1)/4} \cdot (\log d)^{-1/4} \cdot \kappa^{-1/2} & \text{if } s^* \geq 2, \\ d^{-1/2} \cdot (\log d)^{-1/2} \cdot \kappa^{-1/2} & \text{if } s^* = 1. \end{cases} \end{aligned} \quad (\text{E.3})$$

933 On the other hand, we have by Proposition E.3 and Eq. (E.2) that

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g, \theta^* \rangle] &\gtrsim \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}, \theta^* \rangle] - |\langle g, \theta^* \rangle - \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}, \theta^* \rangle]| \\ &\geq |\rho| \gamma (|\rho| \gamma + d^{-1/2})^{s^*-2} - d^{-9s^*} \\ &\geq d^{-(2s^*-1)/4} \cdot (\log d)^{1/4}, \end{aligned}$$

934 whenever $|\rho| \gamma = \Omega(d^{-3/4})$. Therefore, when κ is sufficiently large, we have the fluctuations to be
 935 strictly bounded by half of the signal strength. Thus, we have

$$|\langle g, \theta^* \rangle| \geq \frac{1}{2} \cdot |\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}], \theta^* \rangle|.$$

936 For the norm of g , we have that

$$\begin{aligned} \|g\|_2 &\leq \left| \|g\|_2 - \|\bar{g}\|_2 \right| + \left| \|\bar{g}\|_2 - \|\mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}]\|_2 \right| \\ &\quad + \left| \|\mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}]\|_2 - \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}, \theta^* \rangle] \right| + \left| \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}, \theta^* \rangle] \right| \\ &\leq d^{-9s^*} + \|\bar{g} - \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}]\|_2 + \left| \|\mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}]\|_2 - \langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}], \theta^* \rangle \right| + |\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}], \theta^* \rangle| \\ &\lesssim \left| \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}, \theta^* \rangle] \right| + (\gamma |\rho| + d^{-1/2})^{s^*} + \sqrt{\frac{(d^{-1} \log d)^{(s^*-1)/2} \cdot d \log d}{n}}, \end{aligned} \quad (\text{E.4})$$

937 where in the second inequality, we apply Eq. (E.2) and the triangular inequality that $\left| \|\bar{g}\| - \right|$
 938 $\|\mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}]\| \leq \|\bar{g} - \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}]\|$. And the last inequality is deduced by combining the result Propo-
 939 sition E.3 and the fact that $d^{-9s^*} \ll d^{-s^*/2}$. For the leading term, it holds by Proposition E.3
 940 that

$$\left| \langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}], \theta^* \rangle \right| \simeq \begin{cases} (|\rho| \gamma + d^{-1/2})^{s^*-1} & \text{if } s^* \text{ is odd,} \\ |\rho| \gamma (|\rho| \gamma + d^{-1/2})^{s^*-2} & \text{if } s^* \text{ is even.} \end{cases} \quad (\text{E.5})$$

941 Recall that the alignment admits the following iterative update rule:

$$|\langle \theta', \theta^* \rangle| = \left| \left\langle \frac{\theta + \eta G}{\|\theta + \eta G\|_2}, \theta^* \right\rangle \right| \geq \frac{|\langle G, \theta^* \rangle| - \eta^{-1} |\langle \theta, \theta^* \rangle|}{1 + \eta^{-1}},$$

942 where $G = g/\|g\|_2$. In the following, we will define $\rho^* = d^{-1/4}(\log d)^{1/4}$ as a critical threshold
 943 before the weak alignment. Specifically, in the phase I of weak alignment, we assume that $|\rho| \leq \rho^*$.
 944 When the training process goes across this critical threshold, the dominant term in $\mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}, \theta^* \rangle] \simeq$
 945 $(\gamma \rho)^{\mathbb{1}\{s^* \text{ is even}\}} (\gamma \rho + d^{-1/2})^{s^*-1 - \mathbb{1}\{s^* \text{ is even}\}}$ becomes $\gamma \rho$ instead of $d^{-1/2}$.

946 **Stage I of weak alignment. Case I: $s^* \neq 1$, s^* is odd.** Combining the results in Eq. (E.3), (E.4)
 947 and (E.5), we conclude that for $s^* \geq 3$ being odd,

$$\begin{aligned} |\langle g, \theta^* \rangle| &\gtrsim (|\rho| \cdot (d \log d)^{-1/4} + d^{-1/2})^{s^*-1}, \\ \|g\|_2 &\lesssim (|\rho| \cdot (d \log d)^{-1/4} + d^{-1/2})^{s^*-1} + d^{-(2s^*-3)/4} \cdot (\log d)^{-1/4} \cdot \kappa^{-1/2}. \end{aligned}$$

948 It is important to note that here “ \gtrsim ” and “ \lesssim ” only hides constants that are independent of d and n .
 949 Combining these two inequalities, we have that

$$\begin{aligned} \frac{|\langle g, \theta^* \rangle|}{\|g\|_2} &\gtrsim \frac{(|\rho| d^{1/4} (\log d)^{-1/4} + 1)^{s^*-1}}{(|\rho| d^{1/4} (\log d)^{-1/4} + 1)^{s^*-1} + d^{1/4} (\log d)^{-1/4} \cdot \kappa^{-1/2}} \\ &= \frac{(|\rho|/\rho^* + 1)^{s^*-1}}{(|\rho|/\rho^* + 1)^{s^*-1} + \kappa^{-1/2}/\rho^*}. \end{aligned}$$

950 Thus, if $|\rho| \leq \rho^*$ and take κ to be a sufficiently large constant, after the first gradient update, the
 951 alignment will grow to at least ρ^* by noting that $|\langle g, \theta^* \rangle|/\|g\|_2 \gtrsim \rho^* \kappa^{1/2}$ and that

$$\begin{aligned} |\rho'| &\gtrsim \frac{\rho^* \kappa^{1/2} - \eta^{-1} |\rho|}{1 + \eta^{-1}} \\ &\geq \rho^* \cdot \frac{\sqrt{\kappa} - \eta^{-1}}{1 + \eta^{-1}} \geq \rho^*. \end{aligned}$$

952 As a summary of Case I(a), with one step of gradient update, the alignment will grow to at least ρ^* if
 953 $|\rho| \leq \rho^*$.

954 **Stage I of weak alignment. Case II: s^* is even.** In the case where s^* is even, we have by the
 955 previous arguments that

$$\begin{aligned} |\langle g, \theta^* \rangle| &\gtrsim |\rho| \cdot (d \log d)^{-1/4} \cdot (|\rho| \cdot (d \log d)^{-1/4} + d^{-1/2})^{s^*-2}, \\ \|g\|_2 &\lesssim |\rho| \cdot (d \log d)^{-1/4} \cdot (|\rho| \cdot (d \log d)^{-1/4} + d^{-1/2})^{s^*-2} \\ &\quad + d^{-(2s^*-3)/4} \cdot (\log d)^{-1/4} \cdot \kappa^{-1/2} \end{aligned}$$

956 Here, we use the following fact that

$$(|\rho| \gamma + d^{-1/2})^{s^*} \lesssim (|\rho| \gamma + d^{-1/2})^{s^*-2} \cdot (\rho^2 \gamma^2 + d^{-1}) \lesssim (|\rho| \gamma + d^{-1/2})^{s^*-2} |\rho| \gamma,$$

957 where the last inequality holds since $|\rho| \geq d^{-1/2}/2$. Thus, we conclude that

$$\begin{aligned} \frac{|\langle g, \theta^* \rangle|}{\|g\|_2} &\gtrsim \frac{|\rho| \cdot (|\rho| d^{1/4} (\log d)^{-1/4} + 1)^{s^*-2}}{|\rho| \cdot (|\rho| d^{1/4} (\log d)^{-1/4} + 1)^{s^*-2} + \kappa^{-1/2}} \\ &= \frac{|\rho| \cdot (|\rho|/\rho^* + 1)^{s^*-2}}{|\rho| \cdot (|\rho|/\rho^* + 1)^{s^*-2} + \kappa^{-1/2}}. \end{aligned} \tag{E.6}$$

958 Note that $\kappa = \Omega(1)$. Hence before the alignment reaches ρ^* , $\kappa^{-1/2}$ will dominate the denominator in
 959 Eq. (E.6), which gives us that $|\langle g, \theta^* \rangle|/\|g\|_2 \gtrsim |\rho| \cdot \sqrt{\kappa}$. Importantly, the “ \gtrsim ” hides constants that
 960 are independent of d and κ . Thus, by taking κ to be a sufficiently large constant, we can conclude that

$$|\rho'| \geq |\rho| \cdot \frac{\sqrt{\kappa} - \eta^{-1}}{1 + \eta^{-1}} \geq 2|\rho|.$$

961 As a summary of Case II(a), before the alignment reaches ρ^* , the alignment will grow exponentially
 962 fast, and this phase takes at most $O(\log(d))$ steps. In the following, we consider the case when
 963 $|\rho| \geq \rho^*$ for both cases I and II.

964 **Stage II of weak alignment. Case I&II combined.** Now we consider the case when $|\rho| \geq \rho^*$
 965 for both cases I and II, i.e., $s^* \geq 2$. For this case, we have $|\rho|/\rho^* + 1 \simeq |\rho|/\rho^*$. Let us define
 966 $r = |\rho|/\rho^* \geq 1$ and $r' = |\rho'|/\rho^*$, and it follows that

$$\frac{|\langle g, \theta^* \rangle|}{\|g\|_2} \gtrsim \frac{r^{s^*-1}}{r^{s^*-1} + \kappa^{-1/2}/\rho^*},$$

967 and consequently:

$$\begin{aligned} r' &\gtrsim \frac{r^{s^*-1} \cdot (r^{s^*-1} \rho^* + \kappa^{-1/2})^{-1} - \eta^{-1} r}{1 + \eta^{-1}} \\ &\geq \frac{(\sqrt{\kappa} \cdot r^{s^*-1}) \wedge \rho^{s^*-1} - \eta^{-1} r}{1 + \eta^{-1}} \end{aligned}$$

968 It can be noted that the maximal ratio $r \leq (\rho^*)^{-1}$, and also $\sqrt{\kappa} \cdot r^{s^*-1} \geq 2\eta^{-1} r$ given that κ is
969 sufficiently large and $r \geq 1$. Thus, we conclude that in this case

$$r' \gtrsim (\sqrt{\kappa} \cdot r^{s^*-1}) \wedge \rho^{s^*-1}.$$

970 For this case, the growth of the alignment will be also at least exponentially fast, until it reaches
971 $\Omega((\rho^*)^{-1})$, i.e., $|\rho| = C$ for some constant C . This phase takes at most $O(\log(d))$ steps.

972 **Strong alignment. Case I&II combined.** We need a more careful analysis for this case in order to
973 achieve strong alignment. When the alignment is on a constant level, we can deduce from its original
974 form that

$$|\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}], \theta^* \rangle| = B(\rho, \{\xi_l\}_{l \in [L]}) \cdot (|\rho| \gamma)^{s^*-1} + E_1,$$

975 where $B(\rho, \{\xi_l\}_{l \in [L]}) = \Omega(1)$ is a constant that depends on ρ and the random perturbations $\{\xi_l\}_{l \in [L]}$
976 and the error term follows that $|E| \leq O(d^{-s^*})$. In the following, we will drop the dependency on ρ
977 and $\{\xi_l\}_{l \in [L]}$ and use B for simplicity. We have by Eq. (E.3) that

$$\begin{aligned} |\langle g, \theta^* \rangle| &\geq B(|\rho| \gamma)^{s^*-1} - O\left(d^{-(2s^*-1)/4} \cdot (\log d)^{-1/4} \cdot \kappa^{-1/2}\right) \\ &= B(|\rho| \gamma)^{s^*-1} - O\left((d \log d)^{-s^*/4} \cdot \gamma^{s^*-1} \cdot \kappa^{-1/2}\right) \\ &= \gamma^{s^*-1} \cdot \left(B|\rho|^{s^*-1} - (d \log d)^{-s^*/4} \cdot O(\kappa^{-1/2})\right), \end{aligned}$$

978 and also

$$\begin{aligned} \|g\|_2 &\leq B(|\rho| \gamma)^{s^*-1} + O\left((|\rho| \gamma)^{s^*} + d^{-(2s^*-3)/4} \cdot (\log d)^{-1/4} \cdot \kappa^{-1/2}\right) \\ &\leq B(|\rho| \gamma)^{s^*-1} + O\left((|\rho| \gamma)^{s^*} + d^{-(s^*-2)/4} \cdot (\log d)^{-s^*/4} \cdot \gamma^{s^*-1} \cdot \kappa^{-1/2}\right) \\ &\leq \gamma^{s^*-1} \cdot (B|\rho|^{s^*-1} + O(d^{-(s^*-2)/4} \cdot (\log d)^{-s^*/4} \cdot \kappa^{-1/2})). \end{aligned}$$

979 Therefore, once the alignment reaches a constant level $|\rho| \geq O(1)$, we have

$$\begin{aligned} |\langle G, \theta^* \rangle| &= \frac{|\langle g, \theta^* \rangle|}{\|g\|_2} \geq \frac{B|\rho|^{s^*-1} - O(d^{-s^*/4})}{B|\rho|^{s^*-1} + O(d^{-1/4}(\log d)^{1/4})} \\ &\geq 1 - O(d^{-1/4}(\log d)^{1/4}) =: 1 - \Delta. \end{aligned}$$

980 Here, $\Delta \simeq d^{-1/4}(\log d)^{1/4}$. Thus, as long as $\eta > 2$, after one step gradient,

$$\begin{aligned} |\rho'|^2 &= \frac{\langle G + \eta^{-1}\theta, \theta^* \rangle^2}{\langle G + \eta^{-1}\theta, \theta^* \rangle^2 + \|P_{\theta^*}^\perp(G + \eta^{-1}\theta)\|_2^2} \\ &\geq \frac{(1 - \Delta - \eta^{-1}|\rho|)^2}{(1 - \Delta - \eta^{-1}|\rho|)^2 + (\sqrt{1 - (1 - \Delta)^2} + \eta^{-1}\sqrt{1 - \rho^2})^2} \\ &\geq \frac{(1 - \eta^{-1} - \Delta)^2}{(1 - \eta^{-1} - \Delta)^2 + \eta^{-2}(1 - \rho^2) + 2\sqrt{2}\Delta + 2\Delta^2} \\ &= \frac{(1 - \eta^{-1})^2}{(1 - \eta^{-1})^2 + \eta^{-2}(1 - \rho^2)} - O(\sqrt{\Delta}). \end{aligned}$$

981 Here, the first equality holds by the Pythagorean theorem, the first inequality holds by the triangle
982 inequality, and in the last line, we separate the major term and the error term that scales with $\sqrt{\Delta}$,

983 where we use the fact that $1 - \eta^{-1} > 1/2$ with $\eta > 2$. In addition, by letting $\tau = \eta^{-2}/(1 - \eta^{-1})^2$,
 984 we have

$$\begin{aligned} 1 - (\rho')^2 &= 1 - \frac{(1 - \eta^{-1})^2}{(1 - \eta^{-1})^2 + \eta^{-2}(1 - \rho^2)} + O(\sqrt{\Delta}) \\ &= \frac{\tau(1 - \rho^2)}{1 + \tau(1 - \rho^2)} + O(\sqrt{\Delta}) \\ &\leq \tau(1 - \rho^2) + O(\sqrt{\Delta}). \end{aligned}$$

985 Therefore, we conclude that as long as $\tau < 1$, i.e., $\eta > 2$, $1 - \rho^2$ will exponentially decrease to
 986 $O((1 - \tau)^{-1} \cdot \sqrt{\Delta})$, and achieves strong alignment in $O((\log \Delta^{-1})/(\log \tau^{-1}))$ steps.

987 **Weak & strong alignment. Case III: $s^* = 1$.** In this case, we conclude from the previous
 988 arguments that regardless of the alignment level, it always holds that

$$|\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}], \theta^* \rangle| = B = O(1),$$

989 which gives us

$$|\langle g, \theta^* \rangle| \geq B - O(d^{-1/2} \cdot (\log d)^{-1/2} \cdot \kappa^{-1/2}),$$

990 and

$$\|g\|_2 \leq B + O(d^{-1/4}(\log d)^{1/4} + (\log d)^{-1/2} \cdot \kappa^{-1/2}),$$

991 where we use the fact that $n = \kappa \cdot d(\log d)^2$. Therefore, we also have

$$|\langle G, \theta^* \rangle| = \frac{|\langle g, \theta^* \rangle|}{\|g\|_2} \geq 1 - O((\log d)^{-1/2} \cdot \kappa^{-1/2}) = 1 - \Delta,$$

992 where in this case, we also have $\Delta \simeq (\log d)^{-1/2}$ just like $s^* = 2$ in the previous case, and the rest
 993 of the proof follows the same arguments as in the previous case for the strong alignment. \square

994 E.3 Proof of Key Results

995 *Proof of Lemma E.1.* We begin with proving the first part of the lemma. For conciseness, we drop the
 996 superscript (t) and simply denote θ_m as the present weight. We have the projection of $w_{m,l}$ onto θ^* as

$$|\langle w_{m,l}, \theta^* \rangle| = \left| \frac{\gamma \langle \theta_m, \theta^* \rangle + \langle \xi_{m,l}, \theta^* \rangle}{\|\gamma \theta_m + \xi_{m,l}\|_2} \right| \leq 2(\gamma |\rho_m| + \epsilon).$$

997 For direction $v_{m,2}$, by definition we have

$$\begin{aligned} |\langle w_{m,l}, v_{m,2} \rangle| &= \left| \frac{\gamma \langle \theta_m, v_{m,2} \rangle + \langle \xi_{m,l}, v_{m,2} \rangle}{\|\gamma \theta_m + \xi_{m,l}\|_2} \right| \\ &\leq 2 \cdot \left(\gamma \left\langle \theta_m, \frac{\theta_m - \rho_m \theta^*}{\|\theta_m - \rho_m \theta^*\|_2} \right\rangle + \epsilon \right) \\ &= 2(\gamma \sqrt{1 - \rho_m^2} + \epsilon). \end{aligned}$$

998 For the remaining directions, we always have

$$|\langle w_{m,l}, v_{m,i} \rangle| = \left| \frac{\gamma \langle \theta_m, v_{m,i} \rangle + \langle \xi_{m,l}, v_{m,i} \rangle}{\|\gamma \theta_m + \xi_{m,l}\|_2} \right| \leq 2 \cdot \epsilon,$$

999 where we use the fact that $\langle \theta_m, v_{m,i} \rangle = 0$ for $i \geq 2$. This completes the proof for the first part.

1000 On the joint nice event $\tilde{\mathcal{E}}_m(\tilde{\epsilon})$, we have that

$$\begin{aligned} |\langle w_{m,l}, w_{m,l'} \rangle| &\leq 4 \cdot \langle \gamma \theta_m + \xi_{m,l}, \gamma \theta_m + \xi_{m,l'} \rangle \\ &\leq 4(\gamma^2 + \tilde{\epsilon} + \gamma \langle \theta_m, \xi_{m,l} \rangle + \gamma \langle \theta_m, \xi_{m,l'} \rangle). \end{aligned}$$

1001 On the other hand, it holds on the event $\mathcal{E}_m(\epsilon)$ that

$$|\langle \theta_m, \xi_{m,l} \rangle| = |\langle \sqrt{1 - \rho^2} \cdot v_{m,2} + \rho v_{m,1}, \xi_{m,l} \rangle| \leq (\sqrt{1 - \rho^2} + |\rho|) \cdot \epsilon \leq 2\epsilon.$$

1002 Therefore, we have for any $l \neq l'$ that

$$|\langle w_{m,l}, w_{m,l'} \rangle| \leq 4(\gamma^2 + 4\epsilon\gamma + \tilde{\epsilon}).$$

1003 \square

1004 *Proof of Proposition E.3.* Invoking Lemma H.2 with the fact that $\|\theta^*\|_2 = 1$, we can decompose
 1005 $\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g], \theta^* \rangle$ as

$$\begin{aligned} \langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g], \theta^* \rangle &= \sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)] \cdot \langle w_l, \theta^* \rangle^{s+1} \\ &\quad + \sum_{s \geq s^*} \frac{\sqrt{s}}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s-1}(y)] \cdot \langle w_l, \theta^* \rangle^{s-1} \\ &= \mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] \cdot \frac{\sqrt{s^*}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^{s^*-1} + R, \end{aligned} \quad (\text{E.7})$$

1006 where all the remainder terms are collected by R , defined as

$$R = \sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \left(\mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)] \langle w_l, \theta^* \rangle + \mathbb{E}_{\mathbb{Q}}[\zeta_{s+1}(y) \cdot \hat{\psi}_s(y)] \right) \cdot \langle w_l, \theta^* \rangle^s.$$

1007 Below we will analyze the scale of each term in Eq. (E.7), and show that the remainder R is negligible
 1008 compared to the first term in Eq. (E.7) with high probability over the randomness of the injected noise
 1009 ξ_1, \dots, ξ_L .

1010 **Analysis for the remainder term R in Eq. (E.7).** To bound $|R|$, we apply the triangle inequality
 1011 with the fact that $|\langle w_l, \theta^* \rangle| \leq 1$ to get

$$|R| \leq \sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}} \left[|\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)| + |\zeta_{s+1}(y) \cdot \hat{\psi}_s(y)| \right] \cdot |\langle w_l, \theta^* \rangle|^s.$$

1012 Since $\mathbb{E}_{\mathbb{Q}}[\zeta_{s+1}(y)^2] \leq 1$ for all $s \geq 0$ by the property of the decomposition of the likelihood ratio,
 1013 we have

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}}[|\zeta_{s+1}(y) \cdot \hat{\psi}_s(y)|] &\leq \mathbb{E}_{\mathbb{Q}}[\zeta_{s+1}(y)^2]^{1/2} \cdot \mathbb{E}_{\mathbb{Q}}[\hat{\psi}_s(y)^2]^{1/2} \\ &\leq \mathbb{E}_{\mathbb{Q}}[\hat{\psi}_s(y)^2]^{1/2} \\ &\leq \sqrt{\sum_{s=0}^{\infty} \mathbb{E}[\hat{\psi}_s(y)^2]} = O(1), \end{aligned}$$

1014 and similarly for $\mathbb{E}_{\mathbb{Q}}[|\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)|]$. It then suffices to bound $\sum_{s \geq s^*} (\sqrt{s+1})/L \cdot$
 1015 $\sum_{l=1}^L |\langle w_l, \theta^* \rangle|^s$. Recall that we restrict ourselves to the following nice event

$$\mathcal{E}(\epsilon) : \left\{ |\langle \xi_l, \theta^* \rangle| < \epsilon, \quad \max_{2 \leq i \leq d} |\langle \xi_l, v_i \rangle| < \epsilon, \quad \forall l \in [L] \right\},$$

1016 where $\{v_1 = \theta^*, v_2 = (\theta - \rho\theta^*)/\sqrt{1-\rho^2}, v_3, \dots, v_d\}$ is an orthonormal basis. Since we assume
 1017 that $\gamma = o(1)$, it follows from Lemma E.1 that $|\langle w_l, \theta^* \rangle| < 1/2$ for all $l \in [L]$ on $\mathcal{E}(\epsilon)$. Consequently,
 1018 it holds on $\mathcal{E}(\epsilon)$ that

$$\begin{aligned} \sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L |\langle w_l, \theta^* \rangle|^s &\leq \sum_{s \geq s^*} \sqrt{s+1} \cdot \left(\frac{1}{2}\right)^{s-s^*} \cdot \frac{1}{L} \sum_{l=1}^L |\langle w_l, \theta^* \rangle|^{s^*} \\ &\lesssim \frac{1}{L} \sum_{l=1}^L |\langle w_l, \theta^* \rangle|^{s^*}, \end{aligned} \quad (\text{E.8})$$

1019 where \lesssim hides a constant that depends on s^* . Now it reduces to upper bound the right hand side in
 1020 Eq. (E.8). To proceed, we define

$$\tilde{w}_l = \begin{cases} w_l & \text{if } \sup_i |\langle w_l, v_i \rangle| < \epsilon; \\ 0 & \text{otherwise.} \end{cases}$$

1021 It can be easily verified that $\tilde{w}_l = w_l$ for any $l \in [L]$ on $\mathcal{E}(\epsilon)$, and $\{\tilde{w}_l\}_{l \in [L]}$ is a sequence of
 1022 independent and bounded random vectors. By Lemma E.1, we have that

$$|\langle \tilde{w}_l, \theta^* \rangle| \leq 2\gamma \vee \epsilon.$$

1023 To find its second moment, we have by definition that

$$\begin{aligned} \mathbb{E}_{\tilde{w}_l} [\langle \tilde{w}_l, \theta^* \rangle^{2s^*}] &= \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{2s^*} \cdot \mathbb{1} \{ \sup_i |\langle w_l, v_i \rangle| \leq \epsilon \}] \\ &\leq \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{2s^*}] \\ &\simeq (|\rho|\gamma + d^{-1/2})^{2s^*}, \end{aligned}$$

1024 where the last line holds by Lemma H.4. Therefore, we can apply the Bernstein's inequality
 1025 (Lemma J.1) to the right hand side of Eq. (E.8) restricted to $\mathcal{E}(\epsilon)$. We deduce from it that there
 1026 exists a event $\mathcal{E}_{1,1}$ with $\Pr(\mathcal{E}_{1,1}) \geq 1 - d^{-c}/(MT)$, and it holds on $\mathcal{E}_{1,1} \cap \mathcal{E}(\epsilon)$ that

$$\begin{aligned} \frac{1}{L} \sum_{l=1}^L |\langle w_l, \theta^* \rangle|^{s^*} &= \frac{1}{L} \sum_{l=1}^L |\langle \tilde{w}_l, \theta^* \rangle|^{s^*} \\ &\lesssim \left(1 + \sqrt{L^{-1} \log d}\right) \cdot (|\rho|\gamma + d^{-1/2})^{s^*} + \frac{(\epsilon \vee \gamma)^{s^*} \log d}{L} \\ &\lesssim (|\rho|\gamma + d^{-1/2})^{s^*} \end{aligned} \tag{E.9}$$

1027 Here we use the fact that M and T are at most polynomial in d and the last line holds since we choose
 1028 $L = \Omega\left(d^{1/2} \cdot (\epsilon \vee \gamma)^{s^*} \cdot \log d \vee \log d\right)$.

1029 **Analysis for the dominant term in Eq. (E.7).** We then consider the major term
 1030 $L^{-1} \sum_{l=1}^L \langle w_l, \theta^* \rangle^{s^*-1}$. By the definition of \tilde{w}_l , we can approximate its expectation as follows:

$$\begin{aligned} \mathbb{E}_{\tilde{w}_l} [\langle \tilde{w}_l, \theta^* \rangle^{s^*-1}] &= \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{s^*-1} \cdot \mathbb{1} \{ \sup_i |\langle w_l, v_i \rangle| \leq \epsilon \}] \\ &\simeq \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{s^*-1}] \pm \Pr\left(\sup_i |\langle w_l, v_i \rangle| > \epsilon\right) \\ &\simeq \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{s^*-1}] \pm \Pr(\mathcal{E}(\epsilon)^c), \end{aligned}$$

1031 where we use the fact that $|\langle w_l, \theta^* \rangle|^{s^*-1} \leq 1$ and the event $\{\sup_i |\langle w_l, v_i \rangle| > \epsilon\} \subset \mathcal{E}(\epsilon)^c$. Again,
 1032 we have by Lemma H.4 that

$$\mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{s^*-1}] \simeq \begin{cases} (|\rho|\gamma + d^{-1/2})^{s^*-1} & \text{if } s^* \text{ is odd;} \\ \rho\gamma(|\rho|\gamma + d^{-1/2})^{s^*-2} & \text{if } s^* \text{ is even.} \end{cases}$$

1033 Similarly, we have for the second moment that

$$\begin{aligned} \mathbb{E}_{\tilde{w}_l} [\langle \tilde{w}_l, \theta^* \rangle^{2(s^*-1)}] &= \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{2(s^*-1)} \cdot \mathbb{1} \{ \sup_i |\langle w_l, v_i \rangle| \leq \epsilon \}] \\ &\leq \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{2(s^*-1)}] \\ &\simeq (|\rho|\gamma + d^{-1/2})^{2(s^*-1)}. \end{aligned}$$

1034 Given the boundedness on $\mathcal{E}(\epsilon)$ and the second moment characterization, the Bernstein's inequality
 1035 (Lemma J.1) implies that there exists $\mathcal{E}_{1,2}$ with $\Pr(\mathcal{E}_{1,2}) \geq 1 - d^{-c}/(MT)$. Furthermore, it holds on
 1036 $\mathcal{E}_{1,2} \cap \mathcal{E}(\epsilon)$ that

$$\begin{aligned} \frac{1}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^{s^*-1} &= \frac{1}{L} \sum_{l=1}^L \langle \tilde{w}_l, \theta^* \rangle^{s^*-1} \\ &= \mathbb{E}_{w_l} [\langle w_l, \theta^* \rangle^{s^*-1}] + E, \end{aligned}$$

1037 where the error term E is absolutely bounded as

$$\begin{aligned} |E| &\lesssim (|\rho|\gamma + d^{-1/2})^{s^*-1} \cdot \sqrt{\frac{\log d}{L}} + \frac{(\epsilon \vee \gamma)^{s^*-1} \log d}{L} + \Pr(\mathcal{E}(\epsilon)^c) \\ &\lesssim (|\rho|\gamma + d^{-1/2})^{s^*} + \Pr(\bar{\mathcal{E}}_m^{(t)}(\epsilon)) \\ &\lesssim (|\rho|\gamma + d^{-1/2})^{s^*}. \end{aligned}$$

1038 Here, the second line holds because $L = \Omega\left(\left((\epsilon \vee \gamma)^{s^* - 1} \cdot d^{s^*/2}\right) \log d \vee d \log d\right)$ and the last line
 1039 holds because $\Pr(\overline{\mathcal{E}}(\epsilon)) \leq d^{-s^*/2}$.

1040 So far, we have obtained that on the event $\mathcal{E}(\epsilon) \cap \mathcal{E}_{1,1} \cap \mathcal{E}_{1,2}$, the following holds:

$$\mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g, \theta^* \rangle] \simeq \mathbb{E}_{w_l}[\langle w_l, \theta^* \rangle^{s^* - 1}] + E + R,$$

1041 where $|E| + |R| \lesssim (|\rho|\gamma + d^{-1/2})^{s^*}$ given our configuration of L and $\Pr(\mathcal{E}(\epsilon)^c)$. On the other hand,
 1042 provided that $|\rho|\gamma \gg d^{-1}$, we have that

$$\begin{aligned} (|\rho|\gamma + d^{-1/2})^{s^*} &= (|\rho|\gamma + d^{-1/2})^{s^* - 2} \cdot ((|\rho|\gamma)^2 + d^{-1} + 2 \cdot |\rho|\gamma \cdot d^{-1/2}) \\ &= (|\rho|\gamma + d^{-1/2})^{s^* - 2} \cdot |\rho|\gamma \cdot (|\rho|\gamma + d^{-1} \cdot (|\rho|\gamma)^{-1} + 2d^{-1/2}) \\ &\ll (|\rho|\gamma + d^{-1/2})^{s^* - 2} \cdot |\rho|\gamma. \end{aligned}$$

1043 Therefore, $\mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g, \theta^* \rangle]$ is always the major term no matter whether s^* is even or odd, and we have
 1044 that

$$\mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g, \theta^* \rangle] \simeq \begin{cases} \gamma\rho \cdot (\gamma|\rho| + d^{-1/2})^{s^* - 2} & \text{if } s^* \text{ is even;} \\ (\gamma|\rho| + d^{-1/2})^{s^* - 1} & \text{if } s^* \text{ is odd.} \end{cases}$$

1045 We now turn to the norm of $\mathbb{E}_{\mathbb{P}_{\theta^*}}[g]$. We have already shown the projection of $\mathbb{E}_{\mathbb{P}_{\theta^*}}[g]$ onto θ^* . Next,
 1046 define $P_{\theta^*}^\perp = I - \theta^* \theta^{*\top}$ as the projection matrix onto the orthogonal complement of the space
 1047 spanned by θ^* . Now, it follows from Eq. (H.6) that

$$\begin{aligned} \|P_{\theta^*}^\perp \mathbb{E}_{\mathbb{P}_{\theta^*}}[g]\|_2 &\leq \sum_{s \geq s^*} |\mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)]| \cdot \left\| \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^s \cdot w_l \right\|_2 \\ &\lesssim \sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L |\langle w_l, \theta^* \rangle|^s. \end{aligned}$$

1048 Here, the first inequality holds by noting that the second term in Eq. (H.6) lies exactly along the
 1049 direction of θ^* and thus does not contribute to the norm of $P_{\theta^*}^\perp \mathbb{E}_{\mathbb{P}_{\theta^*}}[g]$, while for the first term in
 1050 Eq. (H.6), we use the triangle inequality and the fact that $\|P_{\theta^*}^\perp v\|_2 \leq \|v\|_2$ for any $v \in \mathbb{R}^d$. In the
 1051 second inequality, we also use the triangle inequality and the fact that $\|w_l\|_2 = 1$ for all $l \in [L]$. Here,
 1052 the “ \lesssim ” hides a constant that depends on the boundedness of $\mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)]$ as we have shown
 1053 in the previous analysis.

1054 Note that the term $\sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L |\langle w_l, \theta^* \rangle|^s$ is already handled in Eq. (E.8) and (E.9) under the
 1055 success of event $\mathcal{E}(\epsilon) \cap \mathcal{E}_{1,1}$, on which we have

$$\|\mathbb{E}_{\mathbb{P}_{\theta^*}}[g]\|_2 - |\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[g], \theta^* \rangle| \lesssim \sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L |\langle w_l, \theta^* \rangle|^s \lesssim (|\rho|\gamma + d^{-1/2})^{s^*}.$$

1056 Setting $\mathcal{E}_1 = \mathcal{E}_{1,1} \cap \mathcal{E}_{1,2}$ gives the desired event.

1057 □

1058 *Proof of Proposition E.4.* The polynomial-tail property allows us to control the fluctuation of the
 1059 gradient estimator g in each direction at the level that is determined by the sample size n and the
 1060 corresponding variance. To this end, we begin with calculating the variance of g along each direction.

1061 **Calculating the second moment.** Given θ and θ^* , recall that we consider the following d orthonor-
 1062 mal directions:

$$\theta^*, v_2, v_3, \dots, v_d,$$

1063 where we set $v_2 = (\theta - \langle \theta, \theta^* \rangle \theta^*) / \|\theta - \langle \theta, \theta^* \rangle \theta^*\|_2$ and v_i for $i \geq 3$ are orthogonal to θ^* and v_2 .
 1064 Our goal is to show that g has small variance on each of these directions.

1065 As each sample (z_i, y_i) is independently drawn from \mathbb{P}_{θ^*} , we just need to consider the variance of

$$g_1 = \frac{1}{L} \sum_{l=1}^L (\psi(y_l, \langle w_l, z_1 \rangle) \cdot z_1 - \widehat{\psi}_1(y_l) \cdot w_l)$$

1066 in the direction of v for $v \in \{\theta^*, v_1, \dots, v_{d-1}\}$ as

$$\begin{aligned} \text{Var}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle] &= \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle^2] - \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle]^2 \\ &\leq \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle^2]. \end{aligned}$$

1067 From this we see that it suffices to bound the second moment of $\langle g_1, v \rangle$, which is given by

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle^2] &\lesssim \frac{1}{L^2} \sum_{l, l'=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}}[\psi(y, \langle w_l, z \rangle) \psi(y, \langle w_{l'}, z \rangle) \langle z, v \rangle^2] + \frac{1}{L^2} \sum_{l, l'=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}}[\widehat{\psi}_1(y)^2 \langle w_l, v \rangle \langle w_{l'}, v \rangle] \\ &= \frac{1}{L^2} \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}} \left[\psi(y, \langle w_l, z \rangle) \psi(y, \langle w_{l'}, z \rangle) \langle z, v \rangle^2 \cdot \left(1 + \sum_{s \geq s^*} \zeta_s(y) h_s(\langle \theta^*, z \rangle) \right) \right] \\ &\quad + \frac{1}{L^2} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}}[\psi(y, \langle w_l, z \rangle) \psi(y, \langle w_l, z \rangle) \langle z, v \rangle^2] + \frac{1}{L^2} \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}}[\widehat{\psi}_1(y)^2 \langle w_l, v \rangle \langle w_{l'}, v \rangle] \\ &\quad + \frac{1}{L^2} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}}[\widehat{\psi}_1(y)^2 \langle w_l, v \rangle^2]. \end{aligned}$$

1068 As $\psi(y, z)$ is quadruple-integrable by Assumption 4.1, the above integral is well-defined. We split
1069 the summation into two parts: $l = l'$ and $l \neq l'$. For $l = l'$, we directly have an $O(L^{-1})$ bound for
1070 each direction $\theta^*, v_2, \dots, v_d$ thanks to the polynomial-like tail property of ψ in Assumption 4.1.

1071 For $l \neq l'$, Lemma E.1 implies that we have on the nice event $\tilde{\mathcal{E}}(\tilde{\epsilon})$ that

$$\begin{aligned} |\langle w_l, w_{l'} \rangle| &\leq 4(\gamma^2 + 2\gamma\epsilon + \tilde{\epsilon}) \\ &\leq 8(\gamma^2 + \tilde{\epsilon}) := \epsilon_2. \end{aligned}$$

1072 Invoking Lemma H.3, for any $v \in \{\theta^*, v_2, \dots, v_d\}$, it holds on $\tilde{\mathcal{E}}(\tilde{\epsilon})$ that

$$\begin{aligned} &\mathbb{E}_{\mathbb{Q}} \left[\psi(y, \langle w_l, z \rangle) \psi(y, \langle w_{l'}, z \rangle) \langle z, v \rangle^2 \cdot \left(1 + \sum_{s \geq s^*} \zeta_s(y) h_s(\langle \theta^*, z \rangle) \right) \right] \\ &\lesssim \epsilon_2^{s^*-1} \cdot \left(1 + \frac{\epsilon_1^2}{\epsilon_2} + \left(\frac{\epsilon_1^2}{\epsilon_2} \right)^{s^*-1} \cdot \epsilon + \mathbf{1}(v \perp \theta^*) \cdot \left(\frac{\epsilon_1^2}{\epsilon_2} \right)^{s^*-2} \cdot \frac{\epsilon_0^2}{\epsilon_2} \cdot (\epsilon_1^2 + \epsilon_1 \cdot \mathbf{1}(s^* \geq 4)) \right) \end{aligned} \tag{E.10}$$

1073 where we also define $\epsilon_1 := \max\{|\langle w_l, \theta^* \rangle|, |\langle w_{l'}, \theta^* \rangle|\}$, $\epsilon_0 := \max\{|\langle w_l, v \rangle|, |\langle w_{l'}, v \rangle|\}$. If the nice
1074 event $\mathcal{E}(\epsilon)$ also holds, on which the following holds for all $l \in [L]$:

$$|\langle \xi_l, \theta^* \rangle| < \epsilon, \quad \max_{2 \leq i \leq d} |\langle \xi_l, v_i \rangle| < \epsilon,$$

1075 then we have by Lemma E.1 that

$$|\langle w_l, \theta^* \rangle| \lesssim \gamma|\rho| + \epsilon, \quad |\langle w_l, v_2 \rangle| \lesssim \sqrt{1 - \rho^2} \gamma + \epsilon, \quad |\langle w_l, v_i \rangle| \lesssim \epsilon, \quad \forall i \geq 3, \quad \forall l \in [L].$$

1076 Consequently, we can set $\epsilon_1 \simeq \gamma|\rho| + \epsilon = o(1)$ and

$$\epsilon_0 \simeq \begin{cases} \gamma|\rho| + \epsilon, & \text{if } v = \theta^*, \\ \gamma\sqrt{1 - \rho^2} + \epsilon, & \text{if } v = v_2, \\ \epsilon, & \text{otherwise.} \end{cases}$$

1077 Therefore, we have the ratio

$$\frac{\epsilon_1^2}{\epsilon_2} \simeq \frac{(\gamma|\rho| + \epsilon)^2}{4(\gamma^2 + \tilde{\epsilon})} \simeq \frac{\gamma^2|\rho|^2 + \epsilon^2}{\gamma^2 + \tilde{\epsilon}}, \quad \frac{\epsilon_0^2}{\epsilon_2} \simeq \begin{cases} \frac{(\gamma|\rho| + \epsilon)^2}{8(\gamma^2 + \tilde{\epsilon})} \simeq \frac{\gamma^2|\rho|^2 + \epsilon^2}{\gamma^2 + \tilde{\epsilon}}, & \text{if } v = \theta^*, \\ \frac{(\gamma\sqrt{1 - \rho^2} + \epsilon)^2}{8(\gamma^2 + \tilde{\epsilon})} \simeq \frac{\gamma^2(1 - \rho^2) + \epsilon^2}{\gamma^2 + \tilde{\epsilon}}, & \text{if } v = v_2, \\ \frac{\epsilon^2}{4(\gamma^2 + \tilde{\epsilon})} \simeq \frac{\epsilon^2}{\gamma^2 + \tilde{\epsilon}}, & \text{otherwise.} \end{cases}$$

1078 Since $\epsilon^2 \leq \tilde{\epsilon}$, we can conclude that $\epsilon_1^2/\epsilon_2 \lesssim 1$ and $\epsilon_0^2/\epsilon_2 \lesssim 1$. Hence, the right-hand side of Eq. (E.10)
 1079 is bounded by $\epsilon_2^{s^*-1} \simeq (\gamma^2 + \tilde{\epsilon})^{s^*-1}$ for all $v \in \{\theta^*, v_2, \dots, v_d\}$.

1080 Similarly, let us consider the term $L^{-2} \cdot \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}}[\hat{\psi}_1(y)^2] \langle w_l, v \rangle \langle w_{l'}, v \rangle$. On the good event $\mathcal{E}(\epsilon)$,
 1081 we have

$$\begin{aligned} \frac{1}{L^2} \cdot \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}}[\hat{\psi}_1(y)^2] \langle w_l, v \rangle \langle w_{l'}, v \rangle &\lesssim \epsilon_0^2 \cdot \mathbb{1}\{s^* \leq 2\} \\ &\lesssim \epsilon_2 \mathbb{1}\{s^* \leq 2\} \\ &\lesssim (\gamma^2 + \tilde{\epsilon})^{s^*-1} \cdot \mathbb{1}\{s^* \leq 2\}. \end{aligned}$$

1082 The first inequality holds because $\hat{\psi}_1(y) = 0$ whenever $s^* \geq 2$ because of Assumption 4.1(b) and the
 1083 second inequality holds due to the condition that $\epsilon^2 \leq \tilde{\epsilon}$

1084 Lastly for all the terms that take a single summation over $l \in [L]$, we have them bounded by $1/L$ as
 1085 each term in the summation can be upper bounded by 1. Combining the results for $l = l'$ and $l \neq l'$,
 1086 we have on the event $\mathcal{E}(\epsilon) \cap \tilde{\mathcal{E}}(\tilde{\epsilon})$ that

$$\text{Var}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle] \leq \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle^2] \lesssim (\gamma^2 + \tilde{\epsilon})^{s^*-1} + \frac{1}{L}, \quad \forall v \in \{\theta^*, v_2, \dots, v_d\}.$$

1087 **Concentration.** The first thing is to control the variation of g in the direction of θ^* , where we need
 1088 to upper bound the $L^r(\mathbb{P}_{\theta^*})$ -norm of $\langle g_1, v \rangle$. To this end, we define $G : (\mathbb{R} \times \mathbb{R}) \times \mathbb{R} \rightarrow \mathbb{R}$ as

$$G(z, y, w) = |\psi(y, \langle w, z \rangle) \cdot \langle z, v \rangle| + |\hat{\psi}_1(y) \cdot \langle w, v \rangle|.$$

1089 Also we define the empirical measure $d\mu(w) = L^{-1} \sum_l \delta(w_l)$, then it holds by integral Minkowski's
 1090 inequality that

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle g_1, v \rangle^r]^{1/r} &\leq \left(\int d\mathbb{P}_{\theta^*}(y, z) \left(\int d\mu(w) \cdot G(z, y, w) \right)^r \right)^{1/r} \\ &\leq \int d\mu(w) \left(\int d\mathbb{P}_{\theta^*}(y, z) \cdot G(z, y, w)^r \right)^{1/r} \\ &\lesssim \frac{1}{L} \sum_l \mathbb{E}_{\mathbb{P}_{\theta^*}}[|\psi(y, \langle w_l, z \rangle) \cdot \langle z, v \rangle|^r]^{1/r} + \frac{1}{L} \sum_l |\langle w_l, v \rangle|. \end{aligned} \quad (\text{E.11})$$

1091 For the second term, we have on $\mathcal{E}(\epsilon)$ that $|\langle w_l, v \rangle| \leq 2\gamma + 2\epsilon \leq 1$. Applying the Cauchy-Schwarz
 1092 inequality, we have for the summand of the first term in Eq. (E.11) that

$$\mathbb{E}_{\mathbb{P}_{\theta^*}}[|\psi(y, \langle w_l, z \rangle) \cdot \langle w_l, z \rangle|^r]^{1/r} \leq \mathbb{E}_{\mathbb{P}_{\theta^*}}[|\psi(y, \langle w_l, z \rangle)|^{2r}]^{1/2r} \cdot \mathbb{E}_{\mathbb{P}_{\theta^*}}[|\langle w_l, z \rangle|^{2r}]^{1/2r}.$$

1093 Note that $\mathbb{E}_{\mathbb{P}_{\theta^*}}[|\langle w_l, z \rangle|^{2r}]^{1/2r} \leq (2r-1)!!^{1/(2r)} \lesssim r^{1/2}$, it suffices to deal with $\mathbb{E}_{\mathbb{P}_{\theta^*}}[|\langle w_l, z \rangle|^{2r}]^{1/2r}$.
 1094 To proceed, we can decompose $\langle w_l, z \rangle$ to components that are correlated and independent with y as

$$\begin{aligned} \langle w_l, z \rangle &= \langle w_l - \langle w_l, \theta^* \rangle \theta^*, z \rangle + \langle w_l, \theta^* \rangle \langle \theta^*, z \rangle \\ &= \sqrt{1 - \langle w_l, \theta^* \rangle^2} \cdot x' + \langle w_l, \theta^* \rangle x' \end{aligned}$$

1095 where $x = \langle \theta^*, z \rangle \sim \mathcal{N}(0, 1)$ is independent to $x' = (1 - \langle w_l, \theta^* \rangle^2)^{-1/2} \cdot \langle w_l - \langle w_l, \theta^* \rangle \theta^*, z \rangle \sim$
 1096 $\mathcal{N}(0, 1)$. Therefore, we define a Gaussian noise operator as $U_\rho \psi(y, x) = \mathbb{E}_{x' \sim \mathcal{N}(0, 1)}[\psi(y, \rho x +$
 1097 $\sqrt{1 - \rho^2} x')]$. And it holds that

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}}[\psi(y, \langle w_l, z \rangle)^{2r}] &= \mathbb{E}_{\mathbb{P}}[U_{\langle \theta^*, w_l \rangle} \psi(y, x)^{2r}] \\ &= \mathbb{E}_{\mathbb{Q}} \left[U_{\langle \theta^*, w_l \rangle} \psi(y, x)^{2r} \cdot \frac{\mathbb{P}(x, y)}{\mathbb{Q}(x, y)} \right] \\ &= \mathbb{E}_{\mathbb{Q}} \left[\psi(y, x)^{2r} \cdot U_{\langle \theta^*, w_l \rangle} \left(\frac{\mathbb{P}(x, y)}{\mathbb{Q}(x, y)} \right) \right] \\ &\leq \left(\mathbb{E}_{\mathbb{Q}}[\psi(y, x)^{4r}] \cdot \mathbb{E}_{\mathbb{Q}} \left[\left(U_{\langle \theta^*, w_l \rangle} \left(\frac{\mathbb{P}(x, y)}{\mathbb{Q}(x, y)} \right) \right)^2 \right] \right)^{1/2}, \end{aligned} \quad (\text{E.12})$$

1098 where the second line follows from the property of the Gaussian noise operator in (B.2). By assumption
 1099 of the tail bound in Assumption 4.1, we have that $\mathbb{E}_{\mathbb{Q}}[\psi(y, x)^{4r}] \leq C_p(4r)^{4C_p r}$. For the second term,
 1100 we have by the Parseval's identity that

$$\mathbb{E}_{\mathbb{Q}} \left[\left(\mathbb{U}_{\langle \theta^*, w_l \rangle} \left(\frac{\mathbb{P}(x, y)}{\mathbb{Q}(x, y)} \right) \right)^2 \right] = 1 + \sum_{s \geq s^*} \langle \theta^*, w_l \rangle^2 \cdot \mathbb{E}_{\mathbb{Q}} [\zeta_s(y)^2] \leq 2.$$

1101 where we use the property that on the good event $\mathcal{E}(\epsilon)$ we have $|\langle w_l, \theta^* \rangle| \leq \gamma|\rho| + \epsilon \leq \gamma + \epsilon < 1/2$
 1102 and also $\mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \leq 1$. And in conclusion, we get for the first term in Eq. (E.11) that

$$\frac{1}{L} \sum_l \mathbb{E}_{\mathbb{P}_{\theta^*}} [|\psi(y, \langle w_l, z \rangle) \cdot \langle w_l, z \rangle|^r]^{1/r} \lesssim r^{C_p+1/2}.$$

1103 Combining everything, we have

$$\mathbb{E}_{\mathbb{P}_{\theta^*}} [|\langle g_1, v \rangle|^r]^{1/r} \lesssim r^{C_p+1/2}, \quad \forall v \in \{\theta^*, v_2, \dots, v_d\}.$$

1104 Thus, by Lemma J.3, there exists a $\{(z_i, y_i)\}_{i \in [n]}$ -measurable event $\mathcal{E}_{2,1}$ with $\Pr(\mathcal{E}_{2,1}) \geq 1 -$
 1105 $d^{-c}/(MT)$ and it holds on \mathcal{E}_2 that

$$\begin{aligned} & |\langle g, \theta^* \rangle - \mathbb{E}_{\mathbb{P}_{\theta^*}} [\langle g, \theta^* \rangle]| \\ & \lesssim \sqrt{\frac{\mathbb{E}_{\mathbb{P}_{\theta^*}} [\langle g_1, \theta^* \rangle^2] \cdot \log(d^c MT)}{n}} + \frac{\log(d^c MT) \cdot \log(d^c MT n)^{C_p+1/2}}{n} \\ & \lesssim \sqrt{\frac{((\gamma^2 + \tilde{\epsilon})^{s^*-1} + L^{-1}) \cdot \log(d)}{n}} + \frac{\log(d)^{C_p+3/2}}{n}, \quad \forall v \in \{\theta^*, v_1, \dots, v_{d-1}\}. \end{aligned} \quad (\text{E.13})$$

1106 where we utilize the fact T, M, n all have polynomial dependency on d . Moreover, since we assume
 1107 that

$$n = \Omega \left(((\gamma^2 + \tilde{\epsilon})^{s^*-1} + L^{-1})^{-1} \cdot \log(d)^{2C_p+1} \right),$$

1108 we have that the first term in Eq. (E.13) dominates, which further implies that

$$|\langle g, \theta^* \rangle - \mathbb{E}_{\mathbb{P}_{\theta^*}} [\langle g, \theta^* \rangle]| \lesssim \sqrt{\frac{((\gamma^2 + \tilde{\epsilon})^{s^*-1} + L^{-1}) \cdot \log(d)}{n}}.$$

1109 Meanwhile, for the ℓ_2 -norm of g , we have by the Jensen's inequality that for any $r \geq 1$,

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}} [\|g_1\|_2^r]^{1/r} &= \left(\mathbb{E}_{\mathbb{P}_{\theta^*}} \left[\sum_{v \in \{\theta^*, v_2, \dots, v_d\}} \langle g_1, v \rangle^2 \right]^{r/2} \right)^{1/r} \\ &\leq \sqrt{d} \cdot \left(\frac{1}{d} \cdot \sum_{v \in \{\theta^*, v_2, \dots, v_d\}} \mathbb{E}_{\mathbb{P}_{\theta^*}} [|\langle g_1, v \rangle|^r] \right)^{1/r} \lesssim \sqrt{d} \cdot (r)^{C_p+1/2}. \end{aligned}$$

1110 This polynomial tail bound enables us to apply Lemma J.3 for the ℓ_2 -norm of g , which implies that
 1111 there exists some event $\mathcal{E}_{2,2}$ with $\Pr(\mathcal{E}_{2,2}) \geq 1 - d^{-c}/(MT)$, and it holds on $\mathcal{E}_{2,2}$ that

$$\|g - \mathbb{E}_{\mathbb{P}_{\theta^*}} [g]\|_2 \lesssim \sqrt{\frac{((\gamma^2 + \tilde{\epsilon})^{s^*-1} + L^{-1}) \cdot d \cdot \log(d)}{n}}.$$

1112 Setting $\mathcal{E}_2 = \mathcal{E}_{2,1} \cap \mathcal{E}_{2,2}$ gives the desired event. This concludes the proof of Proposition E.4. \square

1113 **F Proof of the Main Theorem for the Sparse Prior**

1114 **F.1 Proof Outline and Preliminaries**

1115 In this section, we provide a detailed proof for Theorem 5.1. We begin with some good events that
 1116 we will work with.

1117 **Signal concentration.** We begin with a good event on which the signal spreads almost evenly
 1118 within its support. Define a series of event:

$$\begin{aligned}\mathcal{E}_{0,r} &:= \{\|\theta\|_r^r \leq C_r \cdot k^{1-r/2}\}; \\ \mathcal{E}_{0,\infty} &:= \{\|\theta^*\|_\infty \leq C_\infty \cdot k^{-1/2} \log(k)^{1/2}\}; \\ \mathcal{E}_{0,\#} &:= \left\{ \sum_{j \in [d]} \mathbb{1}\{|\theta_j^*| \geq \frac{1}{\sqrt{2k}}\} \geq \frac{k}{4} \right\}.\end{aligned}$$

1119 The following lemma guarantees that the all the events above hold with high probability.

1120 **Lemma F.1** (Good signal). *Suppose that k is sufficiently large such that $k/\log(k) \geq 32c(r \vee 1)$,
 1121 $k/\log^{r+2} k \geq \sqrt{2c} + 2$, then it holds that*

$$\begin{aligned}\Pr(\mathcal{E}_{0,r}) \wedge \Pr(\mathcal{E}_{0,\infty}) &\geq 1 - O(k^{-c_{0,1}}); \\ \Pr(\mathcal{E}_{0,\#}) &\geq 1 - O(\exp\{-c_{0,2}k\}),\end{aligned}$$

1122 for some constants $c_{0,1}, c_{0,2} > 0$.

1123 *Proof of Lemma F.1.* See Appendix F.5. □

1124 For fixed s^* , we collect all the indices r such that the corresponding nice event $\mathcal{E}_{0,r}$ will be involved
 1125 in the coming analysis. Define $S(s^*) = \{s^* - 1, s^* - \mathbb{1}\{s^* \text{ odd}\}, 2s^*, 4s^*\}$. And we will stick to
 1126 the following high probability event

$$\mathcal{E}_0 := \mathcal{E}_{0,\infty} \cap \mathcal{E}_{0,\#} \cap \left(\bigcap_{r \in S(s^*)} \mathcal{E}_{0,r}\right).$$

1127 With Lemma F.1, we have that $\Pr(\mathcal{E}_0) \geq 1 - O(k^{-c_0})$ for some constant $c_0 > 0$.

1128 **Preparation for characterizing one-step gradient.** Following the same manner as the proof for
 1129 the non-sparse case, we first characterize the alignment of the gradient step (without adversarial error
 1130 term $\text{err}_{m,l,i}^{(t)}$). We begin with the definition of a minimal setup, that collects all the essential elements
 1131 to form the one-step gradient. The following definition is the sparse analogue of Definition E.2.
 1132 Apart from the method of generating the noise, in the sparse case, we will analyze the gradient in a
 1133 coordinate-wise manner to adapt to the sparse structure.

1134 **Definition F.2.** Fix k -sparse vectors $\theta, \theta^* \subset \mathbb{S}^{d-1}$ with $\phi = \text{supp}(\theta)$ and $\phi^* = \text{supp}(\theta^*)$. Let
 1135 $\rho = \langle \theta, \theta^* \rangle$. Suppose that a single batch of data $\{(z_i, y_i)\}_{i \in [n]}$ is i.i.d. generated from \mathbb{P}_{θ^*} . We
 1136 fix the index m as the current neuron. We first sample $\phi_{m,1}, \phi_{m,2}, \dots, \phi_{m,L} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathcal{S}_{k,m})$, i.e.,
 1137 uniform distribution over all k -sparse supports with m -th index always included. Given these random
 1138 supports, we sample independent noises $\xi_{m,l} \sim \text{Unif}(\mathbb{S}^{k-1}(\phi_{m,l}))$ for $l \in [L]$. Now, for each $l \in [L]$
 1139 and $\gamma = o(1)$, we define $w_{m,l} = (\gamma\theta + \xi_{m,l})/\|\gamma\theta + \xi_{m,l}\|_2$. Then our target is

$$\bar{g}_m = \frac{1}{nL} \sum_{i=1}^n \sum_{l=1}^L (\psi(y_i, \langle w_{m,l}, z_i \rangle) \cdot z_i - \hat{\psi}_1(y_i) \cdot w_{m,l}). \quad (\text{F.1})$$

1140 The associated good event is defined as

$$\mathcal{E}_m(\epsilon) = \left\{ \sup_{l,j} |\langle \xi_{m,l}, e_j \rangle| \leq \epsilon \right\}; \quad (\text{F.2})$$

$$\tilde{\mathcal{E}}_m = \left\{ \max_{l \neq l'} \left\{ \sup_l |\phi_{m,l} \cap \phi_{m,l'}|, \sup_l |\phi_{m,l} \cap \phi^*|, \sup_l |\phi_{m,l} \cap \phi| \right\} \leq \log k \right\}. \quad (\text{F.3})$$

1141 **Almost orthogonality.** Recall that our perturbation noise $\xi_{m,l}$ is sampled from $\text{Unif}(\mathbb{S}^{k-1}(\phi_{m,l}))$,
 1142 which is approximately isotropic. One can presume that each $\xi_{m,l}$ is evenly distributed among different
 1143 coordinates. Additionally, we expect that $(\xi_{m,l}, \phi_{m,l})$ and $(\xi_{m,l'}, \phi_{m,l'})$ should have a negligible
 1144 overlap. These two qualitative properties, which can help simplify the analysis, are captured by
 1145 Eq. (F.2) and Eq. (F.3) in Definition F.2. The following lemma characterizes the property of the
 1146 perturbed weights $w_{m,l}$ on the nice event $\mathcal{E}_m(\epsilon) \cap \tilde{\mathcal{E}}$.

1147 **Lemma F.3** (Polarized weight on nice event, sparse case). *Consider the setting in Definition F.2 with*
 1148 $\gamma < 1/2$. *Suppose that the nice event $\mathcal{E}_m(\epsilon) \cap \tilde{\mathcal{E}}_m$ holds and $\|\theta^*\|_\infty \leq 1/\log k$, then we have that*

$$\begin{aligned} \sup_{l,j} |\langle w_{m,l}, e_j \rangle| &\leq 2(\gamma|\theta_j| + \epsilon); \\ \sup_l |\langle w_{m,l}, \theta^* \rangle| &\leq 2(\gamma|\rho| + \epsilon). \end{aligned}$$

1149 *Additionally, we have that*

$$\sup_{l \neq l'} |\langle w_{m,l}, w_{m,l'} \rangle| \leq 4(\gamma^2 + \epsilon^2 \log k).$$

1150 *Proof of Lemma F.3.* See Appendix F.5. □

1151 This lemma, serving as the counterpart of Lemma E.1, controls the behavior of the perturbed weight
 1152 $w_{m,l}$ in its coordinates and alignment with θ^* . Additionally, they are approximately orthogonal with
 1153 each other, which allows for good characterization to the second moment of the gradient.

1154 One may notice that the definition of $\tilde{\mathcal{E}}$ differs from the non-sparse case, where we explicitly bound
 1155 the correlation between different $\xi_{m,l}$. This is the benefit of the sparse structure, as two randomly
 1156 sampled k -sparse supports are naturally of low overlap.

1157 Before delving into the component-wise analysis, we begin with a proposition that will be frequently
 1158 used to calculate the average contribution of each term in the gradient. This proposition serves as the
 1159 counterpart of Lemma H.4 in the sparse case.

1160 **Proposition F.4.** *Suppose that the polarization level $\gamma = o(1)$ and the noise $(\xi_{m,l}, \phi_{m,l}) \sim$
 1161 $\text{Unif}(\mathbb{S}^{k-1}(\phi_{m,l})) \otimes \text{Unif}(\mathcal{S}_{k,m})$. Let $\rho = \langle \theta, \theta^* \rangle$ where θ is the polarized direction in $w_{m,l}$. Assume
 1162 that the event $\mathcal{E}_{0,s-1\{s \text{ odd}\}}$ holds. Then we have that*

$$\mathbb{E}_{w_{m,l}} [\langle w_{m,l}, \theta^* \rangle^s] \simeq \begin{cases} \gamma\rho \cdot (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{s-1})^{s-1} & \text{if } s \text{ odd}; \\ (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_s)^s & \text{if } s \text{ even}, \end{cases}$$

1163 where $\delta_s = (k^2/d)^{1/s} = o(1)$ for any $s = O(1)$.

1164 *Proof of Proposition F.4.* See Appendix F.5. □

1165 Recall that we define the good event over the signal as $\mathcal{E}_0 = \mathcal{E}_{0,\infty} \cap \mathcal{E}_{0,\#} \cap \bigcap_{r \in S(s^*)} \mathcal{E}_{0,r}$, where
 1166 $S(s^*) = \{s^* - 1, s^* - 1\{s^* \text{ odd}\}, 2s^*, 4s^*\}$. This definition facilitates our deferred analysis where
 1167 we need to control multiple moments of different orders and we collect all the necessary good events
 1168 in \mathcal{E}_0 in the first place.

1169 F.2 Properties of the Gradient Step

1170 In this section, we preview some properties regarding the gradient defined in Eq. (F.1). The following
 1171 proposition deals with the first moment.

1172 **Proposition F.5** (First-order moment of the gradient). *Suppose that θ^* is fixed such that the nice*
 1173 *event \mathcal{E}_0 holds. Under Definition F.2, we choose $\gamma \leq \epsilon = o(1)$ such that $\Pr(\mathcal{E}_m(\epsilon)^c) \leq O(k^{-s^*})$*
 1174 *and*

$$L = \Omega\left(\log(d) \cdot (k \vee (\epsilon^{s^*-1} \cdot k^{s^*+1}))\right).$$

1175 *Then there exists a $\{\xi_{m,l}\}_{l \in [L]}$ -measurable event $\mathcal{E}_{m,1}$ with $\Pr(\mathcal{E}_{m,1}) \geq 1 - O(k^{-c_1})$ for some*
 1176 *constant $c_1 > 0$, such that on $\mathcal{E}_{m,1} \cap \mathcal{E}_m(\epsilon) \cap \tilde{\mathcal{E}}_m$, it holds for any $j \in [d]$ that*

$$\langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}_m], e_j \rangle = \mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] \cdot \sqrt{s^*} \cdot \mathbb{E}_{\tilde{w}_{m,l}}[\langle \tilde{w}_{m,l}, \theta^* \rangle^{s^*-1}] \theta_j^* + R_{m,j},$$

1177 *where the expectation*

$$\mathbb{E}_{\tilde{w}_{m,l}}[\langle \tilde{w}_{m,l}, \theta^* \rangle^{s^*-1}] \simeq \begin{cases} \gamma\rho \cdot (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{s^*-2})^{s^*-2} \cdot \theta_j^* & \text{if } s^* \text{ is even}; \\ (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{s^*-1})^{s^*-1} \cdot \theta_j^* & \text{if } s^* \text{ is odd}, \end{cases}$$

1178 and $R_{m,j}$ is the remainder that can be bounded by

$$|R_{m,j}| \lesssim \left((k^{-1} \vee (\gamma|\rho| + k^{-1/2}|\theta_m^*|)) \cdot (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+)^{s^*-1} + k^{-s^*} \right) \cdot |\theta_j^*| \\ + (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+)^{s^*} \cdot (\gamma|\theta_j| + k^{-1/2} \cdot (k/d)^{\mathbb{1}\{j \neq m\}/2}) + k^{-(s^*+1)}.$$

1179 *Proof of Proposition F.5.* See Appendix F.4. \square

1180 The statement of this proposition clarifies the leading term explicitly, which enables us to track the
1181 leading term in the strong alignment more precisely. To complete this section, we provide a proposition
1182 that characterizes the fluctuation of the gradient, serving as the counterpart of Proposition E.4.

1183 **Proposition F.6** (Fluctuation of mini-batch gradient). *Under the simplified setting Definition F.2*
1184 *where ψ follows Assumption 4.1. Additionally, suppose that the sample size*

$$n = \Omega\left(((\gamma^2 + \epsilon^2 \log k)^{s^*-1} + L^{-1})^{-1} \cdot \log(d)^{2C_p+2} \right),$$

1185 *where C_p is the order of the polynomial tail in Assumption 4.1(c). Then there exists a $\{(z_i, y_i)\}_{i \in [n]}$ -*
1186 *measurable event $\mathcal{E}_{m,2}$ with $\Pr(\mathcal{E}_{m,2}) \geq 1 - O(d^{-(c+1)}/T)$, such that on $\mathcal{E}_{m,2} \cap \mathcal{E}_m(\epsilon) \cap \tilde{\mathcal{E}}_m$, it*
1187 *holds that*

$$|\langle \bar{g}_m, e_j \rangle - \langle \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}_m], e_j \rangle| \leq \sqrt{\frac{((\gamma^2 + \epsilon^2 \log k)^{s^*-1} + L^{-1}) \cdot \log(d)}{n}},$$

1188 *for any $v \in \{e_1, e_2, \dots, e_d\}$*

1189 *Proof of Proposition F.6.* See Appendix F.4. \square

1190 With these propositions, we have completed the preparation for the analysis of the gradient step and
1191 are ready to move on to the proof of the main theorem.

1192 F.3 Proof of the Main Theorem

1193 *Proof of Theorem 5.1.*

1194 **Preparations.** We first clarify the final good event that we will use throughout the proof. We first
1195 fix $\epsilon = k^{-1/2} \cdot \log k$, then it holds by Lemma J.6 that for each m and t , we have

$$\Pr\left(\mathcal{E}_m^{(t)}(\epsilon)^c\right) \leq Ld \cdot O(\exp\{-ck\} + k^{-\log k/4}).$$

1196 Since L and d are at most polynomials in k , we see that for sufficiently large k , it holds that
1197 $\Pr\left(\mathcal{E}_m^{(t)}(\epsilon)^c\right) \leq k^{-s^*}$. Additionally, we see that $\gamma = k^{-1/2}$ is fixed and our parameter config-
1198 uration

$$n = \Omega\left((k \log^3 k)^{s^*} \cdot \log d\right), \quad L = \Omega(k^{(s^*+3)/2} \cdot \log(k)^{s^*-1})$$

1199 are clearly compatible with the conditions in Proposition F.5 and Proposition F.6. At the t -th step, the
1200 mini-batch $\{(z_i^{(t)}, y_i^{(t)})\}_{i \in [n]}$, $\theta = \theta_m^{(t)}$ and the error-free gradient

$$\bar{g}_m^{(t)} = \frac{1}{nL} \sum_{l=1}^L \sum_{i=1}^n (\psi(y_i^{(t)}, \langle w_{m,l}^{(t)}, z_i^{(t)} \rangle)) \cdot z_i^{(t)} - \hat{\psi}_1(y_i^{(t)}) \cdot w_{m,l}^{(t)} \quad (\text{F.4})$$

1201 together form an instance of Definition F.2, for which we can find the event $\mathcal{E}_{m,1}^{(t)}$ and $\mathcal{E}_{m,2}^{(t)}$, both with
1202 probability at least $1 - O(d^{-c-1}/T)$, such that on $\mathcal{E}_{m,1}^{(t)} \cap \mathcal{E}_{m,2}^{(t)}$ the results in Proposition F.5 and
1203 Proposition F.6 hold.

1204 The gradient we use in the algorithm differs from Eq. (F.4) by the error term $\text{err}_{m,l,i}^{(t)} \cdot z_i^{(t)}$. To control
1205 this difference, we define

$$\mathcal{E}_{m,3}^{(t)} = \left\{ \sup_i \|z_i^{(t)}\| \leq \sqrt{d} \right\}.$$

1206 Given our specification on $\sup_{m,l,i,t} \text{err}_{m,l,i}^{(t)}$, it holds on $\mathcal{E}_{m,3}^{(t)}$ that for any $v \in \mathbb{S}^{d-1}$:

$$\begin{aligned} \|\|g_m^{(t)}\|_2 + \|\bar{g}_m^{(t)}\|_2\| \vee |\langle g_m^{(t)}, v \rangle - \langle \bar{g}_m^{(t)}, v \rangle| &\leq \|g_m^{(t)} - \bar{g}_m^{(t)}\|_2 \\ &\leq \sup_i \|z_i^{(t)}\| \cdot \sup_l \|\text{err}_{m,l,i}^{(t)}\|_2 \\ &\leq d^{-9s^*}. \end{aligned} \quad (\text{F.5})$$

1207 Our final event is fixed to be

$$\mathcal{E} = \mathcal{E}_0 \cap \bigcap_{m \in [d]} \bigcap_{t=1}^T \left(\mathcal{E}_m^{(t)}(\epsilon) \cap \tilde{\mathcal{E}}_m \cap \mathcal{E}_{m,1}^{(t)} \cap \mathcal{E}_{m,2}^{(t)} \cap \mathcal{E}_{m,3}^{(t)} \right).$$

1208 With union bound, we have that $\Pr(\mathcal{E}) \geq 1 - O(d^{-c})$ for some constant $c > 0$.

1209 To avoid confusion, we denote for each m that $\check{g}_m^{(t)} = \mathbb{E}_{\mathbb{P}_{\theta^*}}[\bar{g}_m^{(t)}]$. Later we will encounter some data
1210 dependent index \hat{m} and using $\check{g}_{\hat{m}}^{(t)}$ avoids the ambiguity of the expectation. With our choice of n and
1211 L , Proposition F.6 guarantees that for any m, t, j , it holds that

$$|\bar{g}_{m,j}^{(t)} - \check{g}_{m,j}^{(t)}| \lesssim k^{-(s^*-1/2)} \cdot \log^{-3/2} k. \quad (\text{F.6})$$

1212 In the sequel, we will drop the superscript t whenever there is no ambiguity. We will frequently
1213 involve $\check{g}_m^{(t)} = P_{\text{Top}_k(g_m^{(t)})}(g_m^{(t)})$.

1214 **Weak alignment.** The proof towards the weak alignment in the initial step comprises three parts. In
1215 the first place, we will show that the index of the gradient we choose \hat{m} guarantees that $|\theta_{\hat{m}}^*| \gtrsim k^{-1/2}$.
1216 Thereby, the corresponding gradient exhibits good alignment towards the signal. Based on this, we
1217 can show that the support we choose $\text{Top}_k(g_{\hat{m}})$ is of considerable quality by successfully identifying
1218 $\phi^{**} = \{j : |\theta_j^*| \geq 1/\sqrt{2k}\}$. Combining these elements, we can show that the gradient $\check{g}_{\hat{m}}$ is
1219 well-aligned with the signal θ^* .

1220 We begin with analyzing the quality of $g_{\hat{m}}$ where $\hat{m} = \arg\max_m \|\check{g}_m\|_2$. With this objective in mind,
1221 we first work on deriving a signal-dependent upper bound for $\|\check{g}_{\hat{m}}\|$. Note that $\rho_m = \langle \theta_m, \theta^* \rangle = |\theta_m^*|$,
1222 where $\theta_m = e_m$ is the initial weight. Applying Proposition F.5, we have for any ϕ , $|\phi| = k$ that

$$\begin{aligned} \sum_{j \in \phi} |\check{g}_{\hat{m},j}|^2 &\lesssim (k^{-1/2} |\theta_{\hat{m}}^*|)^2 \mathbb{1}_{\{s^* \text{ even}\}} \cdot (k^{-1/2} |\theta_{\hat{m}}^*| + k^{-1} \delta_+)^{2(s^*-1-\mathbb{1}_{\{s^* \text{ even}\}})} \cdot \sum_{j \in \phi} \theta_j^{*2} \\ &\quad + (k^{-2} \vee (k^{-1/2} |\theta_{\hat{m}}^*|)^2) \cdot (k^{-1/2} |\theta_{\hat{m}}^*| + k^{-1} \delta_+)^{2(s^*-1)} + k^{-2s^*} \sum_{j \in \phi} \theta_j^{*2} \\ &\quad + (k^{-1/2} |\theta_{\hat{m}}^*| + k^{-1} \delta_+)^{2s^*} \cdot k^{-1} \sum_{j \in \phi} \left(\theta_j^2 + \cdot (k/d)^{\mathbb{1}_{\{j \neq m\}}} \right) + k^{-(2s^*+1)} \end{aligned}$$

1223 Note that we have $\|\theta^*\|_\infty \lesssim k^{-1/2} \cdot \log k$, it holds that

$$k^{-2} \vee (k^{-1} \cdot |\theta_{\hat{m}}^*|^2) \cdot (k^{-1/2} |\theta_{\hat{m}}^*| + k^{-1} \delta_+)^{2(s^*-1)} \lesssim k^{-2s^*} \cdot \log(k)^{2s^*}.$$

1224 For any ϕ such that $|\phi| = k$, we have $\sum_{j \in \phi} \theta_j^2 \leq 1$ for any θ such that $\|\theta\|_0 = k$. Therefore, we can
1225 further upper bound this quantity by

$$\begin{aligned} \sup_{|\phi|=k} \sum_{j \in \phi} |\check{g}_{\hat{m},j}|^2 &\lesssim (k^{-1/2} |\theta_{\hat{m}}^*| + k^{-1} \delta_+)^{2s^*-2} + k^{-2s^*} \cdot \log(k)^{2s^*} \\ &\quad + (k^{-1/2} |\theta_{\hat{m}}^*| + k^{-1} \delta_+)^{2s^*} \cdot k^{-1} \\ &\lesssim (k^{-1/2} |\theta_{\hat{m}}^*| + k^{-1} \delta_+)^{2s^*-2} \\ &\quad + k^{-2s^*} \cdot \log^{s^*/2} k + k^{-2s^*} \cdot \log(k)^{2s^*}. \end{aligned} \quad (\text{F.7})$$

1226 Now we combined Eq. (F.7), Eq. (F.6) and Eq. (F.5) to conclude that

$$\begin{aligned} \sup_{|\phi|=k} \sum_{j \in \phi} |g_{\hat{m},j}|^2 &\lesssim k \cdot \sup_j |g_{\hat{m},j} - \bar{g}_{\hat{m},j}|^2 + k \cdot \sup_j |\bar{g}_{\hat{m},j} - \check{g}_{\hat{m},j}|^2 + \sup_{|\phi|=k} \sum_{j \in \phi} |\check{g}_{\hat{m},j}|^2 \\ &\lesssim k^{-(s^*-1)} \cdot (|\theta_{\hat{m}}^*| + k^{-1/2} \delta_+)^{2s^*-2} + k^{-2s^*} \cdot \log(k)^{2s^*}. \end{aligned}$$

1227 By definition of \tilde{g}_m , we can further conclude that

$$\begin{aligned} \|\tilde{g}_{\widehat{m}}\|_2^2 &= \sup_{|\phi|=k} \sum_{j \in \phi} |\langle g_{\widehat{m}}, e_j \rangle|^2 \\ &\lesssim k^{-(s^*-1)} \cdot |\theta_{\widehat{m}}^*|^{2(s^*-1)} + o(k^{-2(s^*-1)}). \end{aligned} \quad (\text{F.8})$$

1228 On the other hand, for $m \in \phi^{**}$, we have $(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+)^r \gtrsim k^{-r}$. Then it holds by
1229 Proposition F.5 that

$$\sum_{j \in \phi^*} |\check{g}_{m,j}|^2 \gtrsim k^{-2(s^*-1)} - \tilde{O}(k^{-2s^*}),$$

1230 and by Eq. (F.6) and (F.5), we have that

$$\|\tilde{g}_m\|_2^2 \gtrsim k^{-2(s^*-1)} - \tilde{O}(k^{-2s^*}). \quad (\text{F.9})$$

1231 Now, combining Eq. (F.9) with Eq. (F.8) by the definition of \widehat{m} , we have that

$$k^{-(s^*-1)}|\theta_{\widehat{m}}^*|^{2(s^*-1)} + o(k^{-2(s^*-1)}) \gtrsim k^{-2(s^*-1)} - \tilde{O}(k^{-2s^*}).$$

1232 We conclude the first step from the last inequality that there exists a global constant $c_1 > 0$ such that
1233 for sufficiently large k :

$$|\theta_{\widehat{m}}^*| \geq \left(c_1 \cdot (1 - o(1)) \cdot k^{-2(s^*-1)} \cdot k^{s^*-1} \right)^{1/2(s^*-1)} \geq c_1' k^{-1/2}. \quad (\text{F.10})$$

1234 Now we move on to the support identification. We have shown that $|\theta_{\widehat{m}}^*| \gtrsim k^{-1/2}$. To establish
1235 $\phi^{**} \subset \widehat{\phi} = \text{Top}_k(g_{\widehat{m}})$, it is sufficient to demonstrate that

$$\sup_{j \notin \phi^*} |g_{\widehat{m},j}| \leq \inf_{j \in \phi^*} |g_{\widehat{m},j}|. \quad (\text{F.11})$$

1236 In the following, we will bound each side separately. Consider $j \notin \phi^*$, we have by Proposition F.5
1237 that

$$|\check{g}_{\widehat{m},j}| \lesssim (k^{-1/2}|\theta_{\widehat{m}}^*| + k^{-1}\delta_+)^{s^*} \cdot (k^{-1/2}|\theta_j| + (k/d)^{-1/2}) + k^{-(s^*+1)}.$$

1238 Combining this upper bound with Eq. (F.10) and Eq. (F.6) gives us that

$$\begin{aligned} |g_{\widehat{m},j}| &\leq |g_{\widehat{m},j} - \bar{g}_{\widehat{m},j}| + |\bar{g}_{\widehat{m},j} - \check{g}_{\widehat{m},j}| + |\check{g}_{\widehat{m},j}| \\ &\lesssim (k^{-1/2} \cdot |\theta_{\widehat{m}}^*| + k^{-1}\delta_+)^{s^*} \cdot (k^{-1/2}|\theta_j| + (k/d)^{-1/2}) + k^{-(s^*+1)} \\ &\quad + d^{-9s^*} + k^{-(s^*-1/2)} \cdot \log^{-3/2} k \\ &\lesssim k^{-(s^*-1/2)} \cdot \log^{-3/2} k + \tilde{O}(k^{-(s^*+1)}). \end{aligned} \quad (\text{F.12})$$

1239 On the other hand, for $j \in \phi^{**}$, we have that

$$\begin{aligned} |\check{g}_{\widehat{m},j}| &\gtrsim |k^{-1/2}\theta_{\widehat{m}}^* + k^{-1}\delta_+|^{s^*-1} \cdot |\theta_j^*| \\ &\gtrsim k^{-(s^*-1/2)} \end{aligned}$$

1240 Similarly, it holds that

$$\begin{aligned} |g_{\widehat{m},j}| &\gtrsim |\check{g}_{\widehat{m},j}| - |g_{\widehat{m},j} - \bar{g}_{\widehat{m},j}| - |\bar{g}_{\widehat{m},j} - \check{g}_{\widehat{m},j}| \\ &\gtrsim k^{-(s^*-1/2)} - o(k^{-(s^*-1/2)}). \end{aligned} \quad (\text{F.13})$$

1241 Comparing Eq. (F.12) and Eq. (F.13), we successfully validate Eq. (F.11) holds, and consequently
1242 $\phi^{**} \subset \widehat{\phi}$.

1243 We complete our proof of weak alignment by analyzing the inner product between $\tilde{g}_{\widehat{m}}$ and θ^* and the
1244 norm of $\tilde{g}_{\widehat{m}}$ respectively. Since $|\theta_{\widehat{m}}^*| \gtrsim k^{-1/2}$, we have that

$$(\gamma|\rho_{\widehat{m}}| + k^{-1/2} \cdot |\theta_{\widehat{m}}^*| + k^{-1}\delta_+)^{s^*-1} \wedge (\gamma|\rho_{\widehat{m}}| + k^{-1/2} \cdot |\theta_{\widehat{m}}^*| + k^{-1}\delta_+)^{s^*-2} \cdot |\gamma\rho_m| \gtrsim k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{s^*-1}$$

1245 For the inner product, we apply Proposition F.5 for \widehat{m} and each $j \in \phi^{**}$ and get that

$$g_{\widehat{m},j} \cdot \theta_j^* \simeq \theta_j^{*2} \cdot k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{s^*-1} \cdot \text{sign}(\theta_{\widehat{m}}^*) \\ - (|R_{\widehat{m},j}| + |g_{\widehat{m},j} - \bar{g}_{\widehat{m},j}| + |\bar{g}_{\widehat{m},j} - \check{g}_{\widehat{m},j}|) \cdot |\theta_j^*|. \quad (\text{F.14})$$

1246 Since $\phi^{**} \subset \widehat{\phi}_1$, we can lower bound the summation of the leading term as

$$\sum_{j \in \widehat{\phi}_1} \theta_j^{*2} \cdot k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{s^*-1} \gtrsim \sum_{j \in \phi^{**}} \theta_j^{*2} \cdot k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{s^*-1} \\ \gtrsim k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{s^*-1} \cdot k^{-1} \cdot |\phi^{**}| \\ \gtrsim k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{s^*-1}, \quad (\text{F.15})$$

1247 where the last line holds by the definition of $\mathcal{E}_{0,\sharp} \subset \mathcal{E}_0$. For the term associated with $R_{\widehat{m},j}$, we have
1248 by the characterization in Proposition F.5 that

$$\sum_{j \in \widehat{\phi}_1} |R_{m,j}| \cdot |\theta_j^*| \lesssim (k^{-s^*} \cdot \log(k)^{s^*} + k^{-s^*}) \cdot \sum_{j \in \phi^{**}} |\theta_j^*|^2 \\ + k^{-s^*} \cdot \log(k)^{s^*} \cdot k^{-1/2} \cdot \left(\sum_{j \in \widehat{\phi}_1} |\theta_j| \cdot |\theta_j^*| + |\theta_j^*| \cdot (k/d)^{\mathbb{1}\{j \neq m\}} \right) \\ + k^{-(s^*+1)} \cdot \sum_{j \in \widehat{\phi}_1} |\theta_j^*|.$$

1249 Applying the Cauchy-Schwarz inequality, we have that

$$\sum_{j \in \widehat{\phi}_1} |\theta_j^*| \cdot (|\theta_j^*| + (k/d)^{\mathbb{1}\{j \neq m\}}) \leq \left(\sum_{j \in [d]} |\theta_j^*|^2 \right)^{1/2} \cdot \left(\left(\sum_{j \in [d]} |\theta_j|^2 \right)^{1/2} + (1 + k^3/d^2)^{1/2} \right) \\ = O(1). \quad (\text{F.16})$$

1250 And therefore

$$\sum_{j \in \widehat{\phi}_1} |R_{m,j}| \cdot |\theta_j^*| \lesssim k^{-s^*} \cdot \log(k)^{s^*} + k^{-s^*} + \widetilde{O}(k^{-(s^*+1/2)}). \quad (\text{F.17})$$

1251 For the rest of the error terms, we have by Eq. (F.6) and Eq. (F.5) that

$$\sum_{j \in \widehat{\phi}_1} (|g_{\widehat{m},j} - \bar{g}_{\widehat{m},j}| + |\bar{g}_{\widehat{m},j} - \check{g}_{\widehat{m},j}|) \cdot |\theta_j^*| \lesssim \sum_{j \in \phi^{**}} |\theta_j^*| \cdot k^{-(s^*-1/2)} \cdot \log(k)^{-3/2} \\ \leq k^{-(s^*-1)} \cdot \log(k)^{-3/2}. \quad (\text{F.18})$$

1252 Combining Eq. (F.14), Eq. (F.15), Eq. (F.17) and Eq. (F.18), we have that

$$|\langle \check{g}_{\widehat{m}}, \theta^* \rangle| \gtrsim k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{s^*-1} - o(k^{-(s^*-1)}).$$

1253 For the norm, we already have in Eq. (F.8) that

$$\|\check{g}_{\widehat{m}}\|_2 \lesssim k^{-(s^*-1)/2} \cdot |\theta_{\widehat{m}}^*|^{(s^*-1)} + o(k^{-(s^*-1)}).$$

1254 Combining last two inequalities, we concludes that

$$\frac{|\langle P_{\widehat{\phi}_1} g_{\widehat{m}}, \theta^* \rangle|}{\|P_{\widehat{\phi}_1} g_{\widehat{m}}\|} \gtrsim \frac{|\theta_{\widehat{m}}^*|^{s^*-1} - o(k^{-(s^*-1)/2})}{|\theta_{\widehat{m}}^*|^{s^*-1} + o(k^{-(s^*-1)/2})} = \Theta(1),$$

1255 given that $|\theta_{\widehat{m}}^*| \geq c'_1 k^{-1/2}$ for some constant $c'_1 > 0$. This concludes the proof of weak alignment.

1256 **Strong Alignment.** Starting from the second step, we have by Algorithm 2 that all the neurons
 1257 share the same weight parameter. Let θ be the weight parameter in any step after the first step and we
 1258 suppose that $\rho = \langle \theta, \theta^* \rangle = \Theta(1)$. From Proposition F.5, we see that now the choice of the gradient
 1259 should not change the quality of the gradient significantly, as $\gamma|\rho| \gg k^{-1/2}|\theta_m^*|$ for any $m \in [d]$.
 1260 Therefore, we start by analyzing the alignment increment using any g_m . This alteration does not
 1261 affect the alignment increment, but substantially simplifies the analysis.

1262 We additionally define a support $\phi^\dagger = \{j \in [d] : |\theta_j^*| \geq k^{-1}\}$. In the following, we fix an arbitrary
 1263 $m \in [d]$. We begin with analyzing the magnitude of $g_{m,j}$ for $j \in \phi^\dagger \setminus \{m\}$. Since $|\rho| = \Omega(1)$, it
 1264 holds that

$$(\gamma|\rho| + k^{-1/2} \cdot |\theta_m^*| + k^{-1}\delta_{s^*-1})^{s^*-1} \simeq (\gamma|\rho| + k^{-1/2} \cdot |\theta_m^*| + k^{-1}\delta_{s^*-2})^{s^*-2} \cdot (\gamma|\rho|) \simeq k^{-(s^*-1)/2}.$$

1265 With triangle inequality, Proposition F.5 indicates that, for $j \in \phi^\dagger$:

$$\begin{aligned} |g_{m,j}| &\geq |\check{g}_{m,j}| - |\bar{g}_{m,j} - \check{g}_{m,j}| - |\bar{g}_{m,j} - g_{m,j}| \\ &\geq k^{-(s^*-1)/2} \cdot |\theta_j^*| - k^{-s^*/2} \cdot |\theta_j^*| \\ &\quad - k^{-(s^*+1)/2} \cdot (|\theta_j| + (k/d)^{\mathbb{1}\{j \neq m\}/2}) - \tilde{O}(k^{-(s^*-1)/2}). \end{aligned}$$

1266 Similarly, we have for $j \notin \phi^*$ that

$$\begin{aligned} |g_{m,j}| &\leq |\check{g}_{m,j}| + |\bar{g}_{m,j} - \check{g}_{m,j}| + |\bar{g}_{m,j} - g_{m,j}| \\ &\lesssim k^{-(s^*+1)/2} \cdot (|\theta_j| + (k/d)^{\mathbb{1}\{j \neq m\}/2}) + \tilde{O}(k^{-(s^*-1)/2}). \end{aligned}$$

1267 Comparing last two inequalities, we have that $|\theta_j| > k^{-1}$ implies that $|g_{m,j}| \geq \max_{j \notin \phi^*} |g_{m,j}|$ for
 1268 sufficiently large k , and therefore

$$\min_{j \in \phi^\dagger} |g_{m,j}| > \max_{j \notin \phi^*} |g_{m,j}|,$$

1269 which means that $\phi^\dagger \subset \hat{\phi}_m = \text{Top}_k(\theta_m)$. Thereby, we have

$$\sum_{j \in \hat{\phi}_m} |\theta_j^*|^2 \geq 1 - \sum_{j \notin \phi^\dagger} |\theta_j^*|^2 \geq 1 - k^{-1}.$$

1270 We then move on to the alignment analysis. Retreat to Eq. (F.23), we define

$$\begin{aligned} \beta_m(\rho, \{\xi_{m,l}\}_{l \in [L]}, \theta^*, \theta) &= \frac{\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] \cdot \sqrt{s^*} \cdot \mathbb{E}_{\tilde{w}_{m,l}}[\langle w_m, l, \theta^* \rangle^{s^*-1}]}{\text{sign}(\rho)^{\mathbb{1}\{s^* \text{ even}\}} \cdot (\gamma|\rho|)^{s^*-1}}; \\ r_{m,j}(\rho, \{\xi_{m,l}\}_{l \in [L]}, \theta^*, \theta) &= \mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}_m, e_j \rangle] - \theta_j^* \cdot (\gamma|\rho|)^{s^*-1} \cdot \text{sign}(\rho)^{\mathbb{1}\{s^* \text{ even}\}} \cdot \beta_m. \end{aligned} \quad (\text{F.19})$$

1271 Then it follows from Proposition F.5 that whenever $|\rho| = \Omega(1)$, we have that $\beta_m > 0$ and

$$\begin{aligned} \beta_m \vee \beta_m^{-1} &< B \\ |r_{m,j}| &\leq r_{m,j}^a + r_{m,j}^b, \end{aligned}$$

1272 where

$$\begin{aligned} r_{m,j}^a &\leq C_a k^{-(s^*+1)/2} \cdot (|\theta_j| + (k/d)^{\mathbb{1}\{j \neq m\}/2}); \\ r_{m,j}^b &\leq C_b \cdot |\theta_j^*| \cdot (\gamma|\rho|)^{s^*} \end{aligned} \quad (\text{F.20})$$

1273 with some global positive constant B, C_a and C_0 , whenever the designated parameters are compatible
 1274 with the definition of our nice event. With this representation, we can deduce by Eq. (F.6) and Eq. (F.5)
 1275 that

$$\begin{aligned} |\langle \check{g}_m, \theta^* \rangle| &= \left| \sum_{j \in \hat{\phi}_m} \langle \check{g}_m, e_j \rangle \cdot \theta_j^* \right| + \sum_{j \in \phi^*} (|\langle g_m, e_j \rangle - \langle \bar{g}_m, e_j \rangle| + |\langle \bar{g}_m, e_j \rangle - \langle \check{g}_m, e_j \rangle|) \cdot |\theta_j^*| \\ &\geq \left| \sum_{j \in \hat{\phi}_m} \langle \check{g}_m, e_j \rangle \cdot \theta_j^* \right| - \sum_{j \in \phi^*} |\theta_j^*| \cdot O(k^{-(s^*-1)/2} \cdot \log^{-3/2} k) \\ &\geq \beta_m \cdot (\gamma|\rho|)^{s^*-1} \cdot \sum_{j \in \hat{\phi}_m} |\theta_j^*|^2 - \sum_{j \in \phi^*} |r_{m,j}| \cdot |\theta_j^*| - O(k^{-(s^*-1)} \cdot \log^{-3/2} k) \\ &\geq (\beta_m \cdot (\gamma|\rho|)^{s^*-1} - C_b \cdot (\gamma|\rho|)^{s^*}) \cdot \sum_{j \in \hat{\phi}_m} |\theta_j^*|^2 - 2C_a \cdot k^{-(s^*+1)/2} - O(k^{-(s^*-1)} \cdot \log^{-3/2} k). \end{aligned} \quad (\text{F.21})$$

1276 where the last inequality holds by Eq. (F.16). To complete the analysis, we need an upper bound for
 1277 $\|\tilde{g}_m\|_2$. Note that by the triangle inequality and Eq. (F.19), we have that

$$\begin{aligned}\|\tilde{g}_m\| &\leq \|P_{\hat{\phi}_m} \check{g}_m\| + \|P_{\hat{\phi}_m} \bar{g}_m - P_{\hat{\phi}_m} \check{g}_m\| + \|P_{\hat{\phi}_m} \bar{g}_m - P_{\hat{\phi}_m} g_m\| \\ &\leq \|\beta_m(\gamma|\rho|)^{s^*-1} \cdot P_{\hat{\phi}_m} \theta^*\| + \|P_{\hat{\phi}_m} r_m^a\| + \|P_{\hat{\phi}_m} r_m^b\| \\ &\quad + O(k^{-(s^*-1)} \cdot \log^{-3/2} k),\end{aligned}$$

1278 where $r_m^a = (r_{m,1}, \dots, r_{m,d}) \in \mathbb{R}^d$ and $r_m^b = (r_{m,1}, \dots, r_{m,d}) \in \mathbb{R}^d$. To proceed, note that by
 1279 Eq. (F.20) and Eq. (F.6), we have that

$$\begin{aligned}\|P_{\hat{\phi}_m} r_m^a\| &\leq C_a k^{-(s^*+1)/2} \cdot (\|\theta\|_2 + \sqrt{1+k^2/d}) \\ &\leq 3C_a k^{-(s^*+1)/2}, \\ \|P_{\hat{\phi}_m} r_m^b\| &\leq C_b \cdot (\gamma|\rho|)^{s^*} \cdot \left(\sum_{j \in \hat{\phi}_m} |\theta_j^*|^2 \right)^{1/2}.\end{aligned}$$

1280 Putting these upper bounds together, we have that

$$\begin{aligned}\|P_{\hat{\phi}_m} g_m\|_2 &\leq (\beta_m \cdot (\gamma|\rho|)^{s^*-1} + C_b(\gamma|\rho|)^{s^*}) \cdot \left(\sum_{j \in \hat{\phi}_m} |\theta_j^*|^2 \right)^{1/2} \\ &\quad + 3C_a k^{-(s^*+1)/2} + O(k^{-(s^*-1)} \cdot \log^{-3/2} k).\end{aligned}\tag{F.22}$$

1281 Note that for any $a_1 \wedge a_2 > b > 0$, it holds that

$$\frac{a_1 - b}{a_2 + b} = \frac{(a_1 - b) \cdot (a_2 - b)}{a_2^2 - b^2} \geq (a_1/a_2 - b/a_2) \cdot (1 - b/a_2).$$

1282 Setting $\Delta = k^{-1} + k^{-(s^*-1)/2} \cdot \log^{-3/2} k \vee k^{-1} = o(1)$, we get by combining Eq. (F.21) and
 1283 Eq. (F.22) that

$$\begin{aligned}\frac{\langle \tilde{g}_m, \theta^* \rangle}{\|\tilde{g}_m\|} &\geq \frac{(\beta_m - C_b \cdot \gamma|\rho|) \cdot \left(\sum_{j \in \hat{\phi}_m} \theta_j^{*2} \right) - 3C_a \cdot k^{-1} - O(k^{-(s^*-1)/2} \cdot \log^{-3/2} k)}{(\beta_m + C_b \cdot \gamma|\rho|) \cdot \left(\sum_{j \in \hat{\phi}_m} \theta_j^{*2} \right)^{1/2} + 3C_a \cdot k^{-1} + O(k^{-(s^*-1)/2} \cdot \log^{-3/2} k)} \\ &\geq (1 - O(\Delta))^{-1} \cdot \left(\frac{1 - C_b \beta_m^{-1} \gamma|\rho|}{1 + C_b \beta_m^{-1} \gamma|\rho|} \cdot \left(\sum_{j \in \hat{\phi}_m} \theta_j^{*2} \right)^{1/2} - O(\Delta) \right) \\ &\geq 1 - C \cdot k^{-1} - O(\Delta)\end{aligned}$$

1284 where the last line holds because $(1 - Ck^{-1})^r \geq 1 - rC \cdot k^{-1}$ for any $C, r > 0$ and sufficiently
 1285 large k . Note that $k^{-1} = O(\Delta)$, we see that

$$\frac{\langle \tilde{g}_m, \theta^* \rangle}{\|\tilde{g}_m\|} \geq 1 - O(\Delta).$$

1286 This concludes the proof of Theorem 5.1.

1287

□

1288 F.4 Proof of the Key Results

1289 *Proof of Proposition F.5.* First, by Lemma H.2, we have for each $j \in [d]$ that

$$\begin{aligned}\langle \mathbb{E}_{\mathbb{P}_{\theta^*}} [\bar{g}_m], e_j \rangle &= \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)] \cdot \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \langle w_{m,l}, \theta^* \rangle^s \cdot \langle w_{m,l}, e_j \rangle \\ &\quad + \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \cdot \hat{\psi}_{s-1}(y)] \cdot \frac{\sqrt{s}}{L} \cdot \sum_{l=1}^L \langle w_{m,l}, \theta^* \rangle^{s-1} \cdot \langle \theta^*, e_j \rangle. \\ &= \mathbb{E}_{\mathbb{Q}} [\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] \cdot \frac{\sqrt{s^*}}{L} \sum_{l=1}^L \langle w_{m,l}, \theta^* \rangle^{s^*-1} \cdot \langle \theta^*, e_j \rangle \\ &\quad + R_1 + R_2,\end{aligned}\tag{F.23}$$

1290 where the remainders R_1, R_2 are defined as

$$R_1 = \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_{s+1}(y) \cdot \hat{\psi}_s(y)] \cdot \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \langle w_{m,l}, \theta^* \rangle^s \cdot \langle \theta^*, e_j \rangle;$$

$$R_2 = \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)] \cdot \frac{\sqrt{s+1}}{L} \cdot \sum_{l=1}^L \langle w_{m,l}, \theta^* \rangle^s \cdot \langle w_{m,l}, e_j \rangle.$$

1291 We also denote the leading signal term as

$$S = \mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y) \cdot \hat{\psi}_{s^*-1}(y)] \cdot \frac{\sqrt{s^*}}{L} \sum_{l=1}^L \langle w_{m,l}, \theta^* \rangle^{s^*-1} \cdot \langle \theta^*, e_j \rangle.$$

1292 By definition R_1 collects the higher order term that aligns with the signal and the R_2 collects all the
1293 terms in the expected gradient that are parallel to $w_{m,l}$. In comparison to the non-sparse case, here we
1294 are analyzing the gradient coordinate-wisely. Therefore, R_1 and R_2 need to be controlled separately.

1295 **Analysis for the dominant term S in Eq. (F.23).** We first define

$$\tilde{w}_{m,l} = w_{m,l} \cdot \mathbb{1}\{\sup_j \langle \xi_{m,l}, e_j \rangle \leq \epsilon\}.$$

1296 By definition of $\tilde{w}_{m,l}$, we have that $\tilde{w}_{m,l}, l \in [L]$ are independent to each other and $\tilde{w}_{m,l} = w_{m,l}$ on
1297 event $\mathcal{E}_m(\epsilon)$. We can approximate the expectation of the $\langle w_{m,l}, \theta^* \rangle$ as

$$\begin{aligned} \mathbb{E}_{\tilde{w}_{m,l}}[\langle \tilde{w}_{m,l}, \theta^* \rangle^{s^*-1}] &= \mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{s^*-1} \cdot \mathbb{1}\{\sup_j \langle \xi_{m,l}, e_j \rangle \leq \epsilon\}] \\ &\simeq \mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{s^*-1}] \pm \Pr(\sup_j \langle \xi_{m,l}, e_j \rangle > \epsilon) \\ &\simeq \mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{s^*-1}] \pm \Pr(\mathcal{E}_m(\epsilon)^c). \end{aligned}$$

1298 Here the last line holds because $\mathcal{E}_m(\epsilon)^c = \cup_l \{\sup_j \langle \xi_{m,l}, e_j \rangle > \epsilon\}$. For the first term, it holds by
1299 Proposition F.4, we have that on $\mathcal{E}_{0, s^*-1 - \mathbb{1}\{s^* \text{ even}\}}$

$$\mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{s^*-1}] \simeq \begin{cases} \gamma\rho \cdot (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{s^*-2})^{s^*-2} & \text{if } s^* \text{ even;} \\ (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{s^*-1})^{s^*-1} & \text{if } s^* \text{ odd.} \end{cases}$$

1300 For the second moment that is involved in the Bernstein's inequality, we have that on $\mathcal{E}_{0, 2s^*-2}$

$$\begin{aligned} \mathbb{E}[\langle \tilde{w}_{m,l}, \theta^* \rangle^{2s^*-2}] &= \mathbb{E}[\langle w_{m,l}, \theta^* \rangle^{2s^*-2} \mathbb{1}\{\sup_j \langle w_{m,l}, e_j \rangle \leq \epsilon\}] \\ &\leq \mathbb{E}[\langle w_{m,l}, \theta^* \rangle^{2s^*-2}] \\ &\simeq (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{2s^*-2})^{2s^*-2} \end{aligned}$$

1301 To proceed, we have by Bernstein's inequality (Lemma J.1) that there exists an event $\mathcal{E}_{m,11}$ with
1302 $\Pr(\mathcal{E}_{m,11}) \geq 1 - O(d^{-c_b,11})$. And it holds on $\mathcal{E}_{m,11} \cap \mathcal{E}_m(\epsilon)$ that

$$\begin{aligned} \frac{1}{L} \sum_l \langle w_{m,l}, \theta^* \rangle^{s^*-1} &= \frac{1}{L} \sum_l \langle \tilde{w}_{m,l}, \theta^* \rangle^{s^*-1} \\ &\simeq \mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{s^*-1}] + E, \end{aligned}$$

1303 where the error term E can be bounded by

$$|E| \leq (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{2s^*-2})^{s^*-1} \cdot \sqrt{\frac{\log(d)}{L}} + \frac{\epsilon^{s^*-1} \log(d)}{L} + \Pr(\mathcal{E}_m^c(\epsilon)).$$

1304 Moreover, the assumption that

$$L \gtrsim \log d \cdot \left(k^2 \vee \left(\epsilon^{s^*-1} \cdot k^{s^*} \right) \right); \quad \Pr(\mathcal{E}_m(\epsilon)^c) \leq k^{-s^*};$$

1305 allows us to simplify the upper bound for E , since

$$|E| \lesssim k^{-1} \cdot (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{2s^*-2})^{s^*-1} + k^{-s^*}. \quad (\text{F.24})$$

1306 In conclusion, we have on $\mathcal{E}_{m,11} \cap \mathcal{E}_m(\epsilon)$ that

$$S \simeq (\mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{s^*-1}] + E) \cdot \langle \theta^*, e_j \rangle.$$

1307 We remain this form for further simplification.

1308 **Analysis for the first remainder R_1 in Eq. (F.23).** For any s, s' , it holds by the property of
 1309 likelihood ratio decomposition that

$$\begin{aligned}\mathbb{E}_{\mathbb{Q}}[|\zeta_s(y) \cdot \widehat{\psi}_{s'}(y)|] &\leq \mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2]^{1/2} \cdot \mathbb{E}_{\mathbb{Q}}[\widehat{\psi}_{s'}(y)^2]^{1/2} \\ &\leq \sqrt{\sum_{s' \geq 0} \mathbb{E}_{\mathbb{Q}}[\widehat{\psi}_{s'}(y)^2]},\end{aligned}$$

1310 and the last quantity is a constant that is independent to s, s' . To bound the summation for $s \geq s^*$, we
 1311 have on by Lemma F.3 that $|\langle w_{m,l}, \theta^* \rangle| \leq \gamma|\rho| + \epsilon < 1/2$ on

$$\begin{aligned}\sum_{s \geq s^*} \frac{\sqrt{s+1}}{L} \sum_{l=1}^L |\langle w_{m,l}, \theta^* \rangle|^s &\lesssim \sum_{s \geq s^*} \sqrt{s+1} \cdot \left(\frac{1}{2}\right)^{s-s^*} \cdot \frac{1}{L} \sum_{l=1}^L |\langle w_{m,l}, \theta^* \rangle|^{s^*} \\ &\lesssim \frac{1}{L} \sum_{l=1}^L |\langle w_{m,l}, \theta^* \rangle|^{s^*}.\end{aligned}\tag{F.25}$$

1312 Now it reduces to bound the right-hand side of Eq. (F.25). Note that on $\mathcal{E}_m(\epsilon)$, $\tilde{w}_{m,l} = w_{m,l}$. We can
 1313 first track the first and second moment of $\langle w_{m,l}, \theta^* \rangle$ as

$$\begin{aligned}\mathbb{E}_{\tilde{w}_{m,l}}[|\langle \tilde{w}_{m,l}, \theta^* \rangle|^{s^*}] &\leq \mathbb{E}_{w_{m,l}}[|\langle w_{m,l}, \theta^* \rangle|^{s^*}] \\ &\leq \mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{2s^*}]^{1/2}.\end{aligned}$$

1314 To bound the last quantity, we see that given $\mathcal{E}_{0,2s^*}$, Proposition F.4 reads

$$\mathbb{E}_{w_{m,l}}[\langle w_{m,l}, \theta^* \rangle^{2s^*}] \leq (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+)^{2s^*}.$$

1315 By Bernstein's inequality, there exists a event $\mathcal{E}_{m,12}$ with $\Pr(\mathcal{E}_{m,12}) \geq 1 - O(d^{-c_b,12})$ such that on
 1316 $\mathcal{E}_{m,12} \cap \mathcal{E}_m(\epsilon)$, it holds that

$$\frac{1}{L} \sum_{l=1}^L |\langle w_{m,l}, \theta^* \rangle|^{s^*} \lesssim \left(1 + \sqrt{\frac{\log d}{L}}\right) \cdot (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+)^{s^*} + \frac{\epsilon^{s^*} \log(d)}{L},$$

1317 Given that $L \gtrsim \log(d) \cdot (k \vee (\epsilon^{s^*} \cdot k^{s^*}))$, it further holds that

$$\frac{1}{L} \sum_{l=1}^L |\langle w_{m,l}, \theta^* \rangle|^{s^*} \lesssim (\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+)^{s^*} + k^{-s^*}.$$

1318 In conclusion, it holds on $\mathcal{E}_{m,12} \cap \mathcal{E}_m(\epsilon)$ that

$$R_1 \lesssim \left((\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+)^{s^*-1} + k^{-s^*}\right) \cdot |\langle \theta^*, e_j \rangle|.$$

1319 **Analysis for the second remainder R_2 in Eq. (F.23).** Similar to Eq. (F.25), we can first upper
 1320 bound R_2 as

$$\begin{aligned}|R_2| &\lesssim \frac{1}{L} \sum_{l=1}^L |\langle w_{m,l}, \theta^* \rangle|^{s^*} \cdot |\langle w_{m,l}, e_j \rangle| \\ &= \frac{1}{L} \sum_{l=1}^L |\langle \tilde{w}_{m,l}, \theta^* \rangle|^{s^*} \cdot |\langle \tilde{w}_{m,l}, e_j \rangle|,\end{aligned}$$

1321 where the last line holds by the definition of $\mathcal{E}_m(\epsilon)$. We decouple the product with the Cauchy-Schwarz
 1322 inequality as follows:

$$\begin{aligned}\mathbb{E}_{\tilde{w}_{m,l}}[|\langle \tilde{w}_{m,l}, \theta^* \rangle|^{s^*} \cdot |\langle \tilde{w}_{m,l}, e_j \rangle|] &\leq \mathbb{E}_{w_{m,l}}[|\langle w_{m,l}, \theta^* \rangle|^{s^*} \cdot |\langle w_{m,l}, e_j \rangle|] \\ &\leq \mathbb{E}_{w_{m,l}}[|\langle w_{m,l}, \theta^* \rangle|^{2s^*}]^{1/2} \cdot \mathbb{E}_{w_{m,l}}[|\langle w_{m,l}, e_j \rangle|^2]^{1/2}\end{aligned}$$

1323 The first term in the upper bound can be tackled with Proposition F.4. For the second term, we have
 1324 that

$$\begin{aligned}\mathbb{E}_{w_{m,l}}[\langle w_{m,l}, e_j \rangle^2] &\lesssim \mathbb{E}[(\langle \xi_{m,l}, e_j \rangle + \gamma \cdot \langle \theta, e_j \rangle)^2] \\ &\lesssim \mathbb{E}[\langle \xi_{m,l}, e_j \rangle^2] + \gamma^2 \theta_j^2 \\ &= \gamma^2 \theta_j^2 + \mathbb{E}[\mathbb{1}\{j \in \phi_{m,l}\} \cdot \xi_{m,l,j}^2] \\ &\lesssim \gamma^2 \theta_j^2 + k^{-1} \cdot (k/d)^{\mathbb{1}\{j \neq m\}}.\end{aligned}$$

1325 Here, the first line holds because $\|\gamma\theta + \xi\|_2 \geq 1/2$. The last line holds by applying Lemma I.5 and
 1326 that $\mathbb{P}(j \in \phi_{m,l}) \leq k/d$ for $j \neq m$. Note that each term in the summation of is bounded by $\epsilon \log(k)$
 1327 up to a constant on $\mathcal{E}_m(\epsilon)$. We have by Bernstein's inequality (Lemma J.1) that, there exists an event
 1328 $\mathcal{E}_{m,13}$ with $\Pr(\mathcal{E}_{m,13}) \geq 1 - O(d^{-c_b,13})$. And it holds on $\mathcal{E}_{m,13} \cap \mathcal{E}_m(\epsilon)$ that

$$\begin{aligned}|R_2| &\lesssim \left(1 + \sqrt{\frac{\log(d)}{L}}\right) \cdot \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+\right)^{s^*} \cdot \left(\gamma|\theta_j| + k^{-1/2}(k/d)^{\mathbb{1}\{j \neq m\}/2}\right) \\ &\quad + \frac{\epsilon^{s^*+1} \log(d)}{L}.\end{aligned}$$

1329 Given that

$$L \gtrsim \log(d) \cdot \left(k \vee (\epsilon^{s^*+1} \cdot k^{s^*+1})\right),$$

1330 we conclude that it holds on $\mathcal{E}_{m,13} \cap \mathcal{E}_m(\epsilon)$ that

$$|R_2| \lesssim \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+\right)^{s^*} \cdot \left(\gamma|\theta_j| + k^{-1/2}(k/d)^{\mathbb{1}\{j \neq m\}/2}\right) + k^{-(s^*+1)}.$$

1331 **Summary of first-order moment.** We now merge previous results to summarize the results for the
 1332 first-order moment. Note that it is sufficient to set

$$L = \Omega\left(\log(d) \cdot \left(k \vee \epsilon^{s^*-1}(k \cdot \log k)^{s^*+1}\right)\right)$$

1333 Define the final event as $\mathcal{E}_{m,1} = \mathcal{E}_{m,11} \cap \mathcal{E}_{m,12} \cap \mathcal{E}_{m,13}$, which is $\{w_{m,l}\}_{l \in [L]}$ measurable. By
 1334 previous analysis, it holds on this event that

$$\begin{aligned}S &\simeq \theta_j^* \cdot (\gamma\rho)^{\mathbb{1}\{s^* \text{ even}\}} \cdot \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{s^*-1-\mathbb{1}\{s^* \text{ even}\}}\right)^{s^*-1-\mathbb{1}\{s^* \text{ even}\}} + \theta_j^* \cdot E; \\ R_1 &\lesssim \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+\right)^{s^*} + k^{-s^*} \cdot |\theta_j^*|; \\ R_2 &\lesssim \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+\right)^{s^*} \cdot \left(\gamma|\theta_j| + k^{-1/2} \cdot (k/d)^{\mathbb{1}\{j \neq m\}/2}\right) + k^{-(s^*+1)}.\end{aligned}$$

1335 Following the error term E in Eq. (F.24), we define $R = R_1 + R_2 + E$, which be bounded by d

$$\begin{aligned}|R| &\lesssim \left(k^{-1} \vee (\gamma|\rho| + k^{-1/2}|\theta_m^*|)\right) \cdot \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+\right)^{s^*-1} + k^{-s^*} \cdot |\theta_j^*| \\ &\quad + \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_+\right)^{s^*} \cdot \left(\gamma|\theta_j| + k^{-1/2} \cdot (k/d)^{\mathbb{1}\{j \neq m\}/2}\right) + k^{-(s^*+1)}.\end{aligned}\tag{F.26}$$

1336 And we summarize the first moment on $\mathcal{E}_{m,1} \cap \mathcal{E}_m(\epsilon)$ as

$$\mathbb{E}_{\mathbb{P}_{\theta^*}}[\langle \bar{g}_m, e_j \rangle] \simeq \theta_j^* \cdot (\gamma\rho)^{\mathbb{1}\{s \text{ even}\}} \cdot \left(\gamma|\rho| + k^{-1/2}|\theta_m^*| + k^{-1}\delta_{s^*-1-\mathbb{1}\{s^* \text{ even}\}}\right)^{s^*-1-\mathbb{1}\{s \text{ even}\}} + R,$$

1337 where R is upper bounded in Eq. (F.26). \square

1338 *Proof of Proposition F.6.* Similar to the proof of Proposition E.4, the proof of this proposition com-
 1339 prises two parts. To begin with, we calculate the variance of each coordinate of \bar{g}_m .

1340 **Second moment calculation.** It suffices to consider the variance of the first sample. To this end,
 1341 we define

$$\bar{g}_{m,1} = \frac{1}{L} \sum_{l=1}^L (\psi(y_1, \langle w_{m,l}, z_1 \rangle) \cdot z_1 - \hat{\psi}_1(y_1) \cdot w_{m,l}).$$

1342 For any $v \in \{e_1, e_2, \dots, e_d\}$, it holds by the definition of $\bar{g}_{m,1}$ that

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}} [\langle \bar{g}_{m,1}, v \rangle^2] &\lesssim \frac{1}{L^2} \sum_{l,l'=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}} [\psi(y, \langle w_l, z \rangle) \psi(y, \langle w_{l'}, z \rangle) \langle z, v \rangle^2] + \frac{1}{L^2} \sum_{l,l'=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}} [\hat{\psi}_1(y)^2 \langle w_l, v \rangle \langle w_{l'}, v \rangle] \\ &= \frac{1}{L^2} \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}} \left[\psi(y, \langle w_l, z \rangle) \psi(y, \langle w_{l'}, z \rangle) \langle z, v \rangle^2 \cdot \left(1 + \sum_{s \geq s^*} \zeta_s(y) h_s(\langle \theta^*, z \rangle) \right) \right] \\ &\quad + \frac{1}{L^2} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}} [\psi(y, \langle w_l, z \rangle) \psi(y, \langle w_l, z \rangle) \langle z, v \rangle^2] + \frac{1}{L^2} \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}} [\hat{\psi}_1(y)^2] \langle w_l, v \rangle \langle w_{l'}, v \rangle \\ &\quad + \frac{1}{L^2} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}} [\hat{\psi}_1(y)^2] \langle w_l, v \rangle^2. \end{aligned} \tag{F.27}$$

1343 In the same manner as the non-sparse case, we can derive a $O(1/L)$ upper bound for the second and
 1344 the last summation, which traverse through all $l = l'$. Recalling Eq. Lemma F.3, we already have that

$$\begin{aligned} \sup_{l,j} |\langle w_{m,l}, e_j \rangle| &\lesssim \gamma + \epsilon; \\ \sup_l |\langle w_{m,l}, \theta^* \rangle| &\lesssim \gamma |\rho| + \epsilon. \\ \sup_{l \neq l'} |\langle w_{m,l}, w_{m,l'} \rangle| &\lesssim \gamma^2 + \epsilon^2 \log(k). \end{aligned}$$

1345 To incorporate with the notations in Lemma H.3, we denote the upper bounds of the $\langle w_{m,l}, e_j \rangle$,
 1346 $\langle w_{m,l}, \theta^* \rangle$ and $\langle w_{m,l}, w_{m,l'} \rangle$ ($l \neq l'$) as

$$\epsilon_0 = \gamma + \epsilon; \quad \epsilon_1 = \gamma |\rho| + \epsilon; \quad \epsilon_2 = \gamma^2 + \epsilon^2 \log(k),$$

1347 respectively. By the virtue of Lemma H.3, the desired expectation is behaving nicely if the ratio
 1348 $(\epsilon_0^2 \vee \epsilon_1^2) / \epsilon_2$ is a constant term. To validate this fact, we note that

$$\frac{\epsilon_0^2}{\epsilon_2} \simeq \frac{\gamma^2 + \epsilon^2}{\gamma^2 + \epsilon^2 \log(k)}, \quad \frac{\epsilon_1^2}{\epsilon_2} \simeq \frac{\gamma^2 \rho^2 + \epsilon^2}{\gamma^2 + \epsilon^2 \log(k)}.$$

1349 Since $\epsilon \ll 1 \ll \log(k)^{1/2}$, we conclude that $\epsilon_0^2 \vee \epsilon_1^2 / \epsilon_2 \lesssim 1$ for sufficiently large k . Therefore, we
 1350 have by Lemma H.3 that

$$\frac{1}{L^2} \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}} \left[\psi(y_1, \langle w_{m,l}, z_1 \rangle) \cdot \psi(y_1, \langle w_{m,l'}, z_1 \rangle) \cdot \langle z_1, v \rangle^2 \cdot \left(1 + \sum_{s=s^*}^{\infty} \zeta_s(y) h_s(\langle \theta^*, z \rangle) \right) \right] \lesssim \epsilon_2^{s^*-1}.$$

1351 On the other hand, we have for the third term in Eq. (F.27) that

$$\begin{aligned} \frac{1}{L^2} \sum_{l \neq l'} \mathbb{E}_{\mathbb{Q}} [\hat{\psi}_1(y)^2] \langle w_{m,l}, v \rangle \cdot \langle w_{m,l'}, v \rangle &\lesssim \sup_l |\langle w_{m,l}, v \rangle|^2 \cdot \mathbf{1}\{s^* \leq 2\} \\ &\lesssim \epsilon_0^2 \cdot \mathbf{1}\{s^* \leq 2\} \\ &\lesssim \epsilon_2 \cdot \mathbf{1}\{s^* \leq 2\}, \end{aligned}$$

1352 where the first line holds by Lemma F.3.

1353 In summary, we have on the event $\mathcal{E}_m(\epsilon) \cap \tilde{\mathcal{E}}_m$ that

$$\sup_{v \in \{e_1, e_2, \dots, e_d\}} n \text{Var}_{\mathbb{P}_{\theta^*}} [\langle g_m, v \rangle] \lesssim \epsilon_2^{s^*-1} + \frac{1}{L} = (\gamma^2 + \epsilon^2 \log(k))^{s^*-1} + \frac{1}{L}.$$

1354 **Concentration .** We now turn to validate the condition of Lemma J.3. For any $v \in \{e_1, e_2, \dots, e_d\}$,
 1355 we set $G(z, y, w) = |\psi(y, \langle w, z \rangle) \cdot \langle z, v \rangle| + |\hat{\psi}_1(y) \cdot \langle w, v \rangle|$ with the domain measure defined as
 1356 $d\mathbb{P}_{\theta^*}(z_1, y_1) \times d\mu(w)$, where $d\mu(w) = L^{-1} \sum_i \delta_{w_{m,i}}$, the integral Minkowski's inequality implies
 1357 that

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}} [|\langle g_{m,1}, v \rangle|^r]^{1/r} &= \left(\int d\mathbb{P}_{\theta^*}(y, z) \left(\int d\mu(w) |G(z, y, w)|^r \right)^{1/r} \right)^{1/r} \\ &\leq \int d\mu(w) \left(\int d\mathbb{P}_{\theta^*}(y, z) |G(z, y, w)|^r \right)^{1/r} \\ &= \frac{1}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}} [|\psi(y_i, \langle w_{m,l}, z_i \rangle) \cdot \langle z_i, v \rangle|^r]^{1/r} + \frac{1}{L} \sum_l |\langle w_{m,l}, v \rangle|. \quad (\text{F.28}) \end{aligned}$$

1358 To proceed, we leverage Cauchy-Schwarz inequality to decouple the average of the product in the
 1359 first term, which reads

$$\frac{1}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}} [|\psi(y_i, \langle w_{m,l}, z_i \rangle) \cdot \langle z_i, v \rangle|^r]^{1/r} \leq \frac{1}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}} [|\psi(y_i, \langle w_{m,l}, z_i \rangle)|^{2r}]^{1/2r} \cdot \mathbb{E}_{\mathbb{P}_{\theta^*}} [|\langle z_i, v \rangle|^{2r}]^{1/2r}.$$

1360 Similar to the proof of Proposition E.4, we have that

$$\mathbb{E}_{\mathbb{P}_{\theta^*}} [|\psi(y, \langle w_{m,l}, z \rangle)|^{2r}] \leq \mathbb{E}_{\mathbb{Q}} [U_{(\theta^*, w_{m,l})} \left(\frac{\mathbb{P}(x, y)}{\mathbb{Q}(x, y)} \right)^2]^{1/2} \cdot \mathbb{E}_{\mathbb{Q}} [\psi(y, x)^{4r}]^{1/2} \lesssim r^{C_p 4r},$$

1361 where the first inequality exactly repeats Eq. (E.12) and the second inequality holds by Assump-
 1362 tion 4.1(c). On the other hand, we have that $\mathbb{E}_{\mathbb{P}_{\theta^*}} [|\langle z_i, v \rangle|^{2r}]^{1/2r} \leq r^{1/2}$. Since the second term in
 1363 Eq. (F.28) is bounded by $O(1)$, we conclude that

$$\mathbb{E}_{\mathbb{P}_{\theta^*}} [|\langle g_{m,1}, v \rangle|^r]^{1/r} \lesssim r^{C_p + 1/2}.$$

1364 Thus, Lemma J.3 implies that there exists a $\{(z_i, y_i)\}_{i \in [n]}$ -measurable event $\mathcal{E}_{m,2}$ with probability at
 1365 least $1 - O(d^{-c-1}/T)$, on which for any $v \in \{e_1, e_2, \dots, e_d\}$, it holds that

$$\begin{aligned} |\langle g_m, v \rangle - \mathbb{E}_{\mathbb{P}_{\theta^*}} [\langle g_m, v \rangle]| &\lesssim \sqrt{\frac{\mathbb{E}_{\mathbb{P}_{\theta^*}} [\langle g_{m,1}, v \rangle^2] \cdot \log(d^{c+1}T)}{n} + \frac{\log(d^{c+1}T) \cdot \log(d^{c+1}Tn)^{C_p + 1/2}}{n}} \\ &\lesssim \sqrt{\frac{((\gamma^2 + \epsilon^2 \log(k))^{s^* - 1} + L^{-1}) \cdot \log(d)}{n} + \frac{\log(d)^{C_p + 3/2}}{n}}, \end{aligned}$$

1366 given that T, n are at most of polynomial rate in d . Since we assume that

$$n = \Omega\left(\left((\gamma^2 + \epsilon^2 \log(k))^{s^* - 1} + L^{-1}\right)^{-1} \cdot \log(d)^{2C_p + 2}\right),$$

1367 the above inequality can be further simplified as

$$|\langle g_m, v \rangle - \mathbb{E}_{\mathbb{P}_{\theta^*}} [\langle g_m, v \rangle]| \lesssim \sqrt{\frac{((\gamma^2 + \epsilon^2 \log(k))^{s^* - 1} + L^{-1}) \cdot \log(d)}{n}}.$$

1368 Additionally, $\mathcal{E}_{m,2}$ is the desired event. This concludes the proof of Proposition F.6. \square

1369 F.5 Proofs for Technical Results in the Sparse Case

1370 *Proof of Lemma F.1.* With slightly abuse of notation, we assume that $\theta^* \sim \text{Unif}(\mathbb{S}^{k-1})$. We first
 1371 consider the event $\mathcal{E}_{0,\infty} = \{\|\theta^*\|_\infty \leq C \cdot k^{-1/2} \log(k)^{1/2}\}$. From the proof of Lemma J.6, we see
 1372 that

$$\mathbb{P}(\|\theta^*\|_\infty \geq t) \leq 2k \cdot \mathbb{P}(\theta_1^* \geq t) \leq 2k \exp(-k/16) + 2k \exp(-t^2 k/4).$$

1373 Take $t = C \cdot k^{-1/2} \log(k)^{1/2}$, we have that the failure probability is upper bounded by
 1374 $2k \exp(-k/16) + 2k^{1-C^2/4}$.

1375 For the r -norm, we leverage the property that $\theta^* \stackrel{d.}{=} Z/\|Z\|_2$ where $Z \sim \mathcal{N}(0, I_k)$. Now,
 1376 $\|Z\|_2^2 = \sum_{i \leq k} Z_i^2$, where $Z_i^2 - \mathbb{E}[Z_i^2] \geq -1$. Applying one-sided Bernstein's inequality with
 1377 failure probability k^{-c_0} , we have that

$$\|Z\|_2^2 \leq k + \sqrt{2c_0 k \log(k)} + c_0 \log(k)/3.$$

1378 On the other hand, note that we have for $c_1 > 2$ that

$$\mathbb{P}(\max_{i \leq k} |Z_i| > \sqrt{2c_1 \log k}) \leq k \cdot \mathbb{P}(|Z_1| > \sqrt{2c_1 \log k}) \leq k \exp\{-c \log k\} = k^{1-c_1},$$

1379 To apply Bernstein's inequality, we note that

$$\begin{aligned} \mathbb{E}[|Z_i|^r \cdot \mathbf{1}\{|Z_1| \leq \sqrt{2c_1 \log k}\}] &\leq \mathbb{E}[|Z_i|^{2r}]^{1/2}, \\ \mathbb{E}[|Z_i|^{2r} \cdot \mathbf{1}\{|Z_1| \leq \sqrt{2c_1 \log k}\}] &\leq \mathbb{E}[|Z_i|^{2r}], \end{aligned}$$

1380 where $\mathbb{E}[Z_i^{2r}] = (2r-1)!!$. Therefore, it holds by truncated Bernstein's inequality that

$$\mathbb{P}\left(\|Z\|_r^r > (k + \sqrt{2c_2 k \log(k)}) \cdot \mathbb{E}[|Z_1|^{2r}]^{1/2} + (\sqrt{2C \log(k)})^r \cdot c_2 \log(k)/3\right) \leq k^{1-c_1} + k^{-c_2}$$

1381 Combining the two bounds, we conclude that with probability $1 - O(k^{1-c_0 \vee c_1 \vee c_2})$, it holds that

$$\|\theta^*\|_r^r \stackrel{d.}{=} \frac{\|Z\|_r^r}{\|Z\|_2^r} \lesssim \frac{k + \sqrt{k \log k} + (\log k)^{1+r/2}}{(k + \sqrt{k \log k} + \log k)^{r/2}} \lesssim k^{1-r/2},$$

1382 with probability at least $1 - O(k^{-c})$ for some constant $c > 0$.

1383 We now move on to consider the event $\mathcal{E}_{0,\#} = \{\sum_{i \leq k} \mathbf{1}\{|\theta_i^*| \geq 1/\sqrt{2k}\} \geq k/4\}$. First, it holds by
 1384 the Hoeffding's inequality that

$$\mathbb{P}\left(\left|\sum_{i \leq k} \mathbf{1}\{Z_i^2 \geq 1/2\} - kp\right| \leq \frac{kp}{2}\right) \geq 1 - 2 \exp(-2p^2k),$$

1385 where $p = \mathbb{P}(Z_1^2 \geq 3/4) > 0.5$. Denote above event as \mathcal{A}_1 . On the other hand, we have by the
 1386 Bernstein's inequality that

$$\mathbb{P}\left(\left|k^{-1} \sum_{i \leq k} Z_i^2 - 1\right| \leq 1/2\right) \geq 1 - 2 \exp\{-k/32\}.$$

1387 Denote above event as \mathcal{A}_2 . Then on the event $\mathcal{A}_1 \cap \mathcal{A}_2$, we have that

$$\begin{aligned} \sum_{i \leq k} \mathbf{1}\left\{\frac{Z_1}{\|Z\|} > \frac{1}{\sqrt{2k}}\right\} &= \sum_{i \leq k} \mathbf{1}\left\{Z_i^2 > \frac{1}{2k} \sum_{i \leq k} Z_i^2\right\} \\ &\geq \sum_{i \leq k} \mathbf{1}\{Z_i^2 > \frac{3}{4}\} \\ &\stackrel{\mathcal{A}_1}{\geq} \frac{kp}{2} > k/4. \end{aligned}$$

1388 In conclusion we have that $\mathbb{P}(\{|i : |\theta_i^*| > 1/\sqrt{2k}|\} > k/4) \geq \mathbb{P}(\mathcal{A}_1 \cap \mathcal{A}_2) \geq 1 - \exp\{-c_3 k\}$ for
 1389 some constant $c_3 > 0$.

1390 □

1391 *Proof of Lemma F.3.* Clearly, it holds that

$$\|\gamma\theta + \xi_{m,l}\|_2 \geq \|\xi_{m,l}\|_2 - \gamma \cdot \|\xi_{m,l}\|_2 \geq 1/2.$$

1392 Bu substituting this lower bound for the denominator, we have for any j, l that

$$\begin{aligned} |\langle w_{m,l}, e_j \rangle| &\leq 2(\gamma|\theta_j| + |\langle \xi_{m,l}, e_j \rangle|) \\ &\leq 2(\gamma|\theta_j| + \epsilon). \end{aligned}$$

1393 The last line holds by the definition of $\mathcal{E}_m(\epsilon)$. On the other hand, we have for any l that

$$\begin{aligned} |\langle w_{m,l}, \theta^* \rangle| &\leq 2(\gamma|\rho| + \sum_{j \in [d]} \xi_{m,l,j} \cdot \theta_j^* \cdot \mathbf{1}\{j \in \phi^* \cap \phi_{m,l}\}) \\ &\leq 2\gamma|\rho| + 2\left(\sum_j \xi_{m,l,j}^2 \cdot \mathbf{1}\{j \in \phi^* \cap \phi_{m,l}\}\right)^{1/2} \cdot \left(\sum_j \theta_j^{*2} \cdot \mathbf{1}\{j \in \phi^* \cap \phi_{m,l}\}\right)^{1/2} \\ &\leq 2\gamma|\rho| + 2 \sup_j |\xi_{m,l,j}| \cdot \|\theta^*\|_\infty \cdot |\phi^* \cap \phi_{m,l}| \end{aligned}$$

1394 To proceed, note that on the event $\tilde{\mathcal{E}}_m \cap \mathcal{E}_m(\epsilon)$, it holds that $|\phi^* \cap \phi_{m,l}| \leq \log k$ and that $\sup_j |\xi_{m,l,j}| \leq$
1395 ϵ . Since we assume that $\|\theta^*\|_\infty \leq 1/\log k$, it holds that

$$|\langle w_{m,l}, \theta^* \rangle| \leq 2(\gamma|\rho| + \epsilon).$$

1396 Now we turn to consider the correlation between $w_{m,l}$ and $w_{m,l'}$.

$$\begin{aligned} |\langle w_{m,l}, w_{m,l'} \rangle| &\leq 2\left(\gamma^2 + \sum_j |\xi_{m,l,j}| \cdot |\xi_{m,l',j}| \cdot \mathbf{1}\{j \in \phi_{m,l} \cap \phi_{m,l'}\}\right) \\ &\quad + \gamma \sum_j |\theta_j| \cdot |\xi_{m,l,j}| \cdot \mathbf{1}\{j \in \phi_{m,l} \cap \text{supp}(\theta)\} \\ &\quad + \gamma \sum_j |\theta_j| \cdot |\xi_{m,l',j}| \cdot \mathbf{1}\{j \in \phi_{m,l'} \cap \text{supp}(\theta)\}. \end{aligned}$$

1397 For the second term, we have with the definition of $\mathcal{E}_m(\epsilon)$ that

$$\begin{aligned} \sum_j |\xi_{m,l,j}| \cdot |\xi_{m,l',j}| \cdot \mathbf{1}\{j \in \phi_{m,l} \cap \phi_{m,l'}\} &\leq \max_{j,l} |\xi_{m,l,j}|^2 \cdot |\phi_{m,l} \cap \phi_{m,l'}| \\ &\leq \epsilon^2 \log k. \end{aligned}$$

1398 For the third term, applying the Cauchy-Schwarz inequality, we have that

$$\gamma \sum_j |\theta_j| \cdot |\xi_{m,l',j}| \cdot \mathbf{1}\{j \in \phi_{m,l} \cap \phi_{m,l'}\} \leq \gamma \|\theta\|_2 \cdot \epsilon \sqrt{\log k} \leq \gamma^2 + \epsilon^2 \log k.$$

1399 Putting them together, we have that

$$|\langle w_{m,l}, w_{m,l'} \rangle| \leq 4(\gamma^2 + \epsilon^2 \log k).$$

1400 This concludes the proof of Lemma F.3. \square

1401 *Proof of Proposition F.4.* For conciseness, we momentarily drop the subscript m, l in the following
1402 analysis. Conditioning on fixed ϕ , we have that

$$\begin{aligned} \mathbb{E}_w[\langle w, \theta^* \rangle^s] &= \mathbb{E}_w[\|\gamma\theta + \xi\|_2^{-s} \cdot (\gamma\langle \theta, \theta^* \rangle + \langle \xi, \theta^* \rangle)^s] \\ &= \mathbb{E}_\phi[\mathbb{E}_w[\|\gamma\theta + \xi\|_2^{-s} \cdot (\gamma\rho + \langle \xi, P_\phi \theta^* \rangle)^s \mid \phi]]. \end{aligned} \quad (\text{F.29})$$

1403 Given the polarization level $\gamma = o(1)$, we see that $\|\gamma\theta + \xi\|_2^{s+1} \simeq 1 \pm o(1)$, and it suffices to evaluate
1404 $\mathbb{E}_w[(\gamma\rho + \langle \xi, P_\phi \theta^* \rangle)^s \mid \phi]$. Without loss of generality, we assume that $1 \in \phi$ and we can translate
1405 $P_\phi \theta^*$ into the first coordinate by the isotropy of ξ over $\mathbb{S}^{k-1}(\phi)$. To this end, we can characterize the
1406 first term as follows:

$$\begin{aligned} \mathbb{E}[(\gamma\rho + \langle \xi, P_\phi \theta^* \rangle)^s \mid \phi] &= \mathbb{E}[(\gamma\rho + \langle \xi, \|P_\phi \theta^*\|_2 \cdot e_1)^s \mid \phi] \\ &= \sum_{r=0}^s \binom{2\lfloor s/2 \rfloor}{r} (\gamma\rho)^{s-r} \cdot \|P_\phi \theta^*\|_2^r \cdot \mathbb{E}[\xi_1^r \mid \phi] \cdot \mathbf{1}\{r \text{ even}\} \\ &\stackrel{(i)}{\simeq} \sum_{r=0}^{\lfloor s/2 \rfloor} \binom{2\lfloor s/2 \rfloor}{2r} (\gamma\rho)^{2\lfloor s/2 \rfloor - 2r} \cdot \|P_\phi \theta^*\|_2^{2r} \cdot k^{-r} \cdot (\gamma\rho)^{\mathbf{1}\{s \text{ odd}\}} \\ &= (\gamma\rho)^{\mathbf{1}\{s \text{ odd}\}} \cdot ((\gamma\rho + k^{-1/2} \|P_\phi \theta^*\|_2)^{2\lfloor s/2 \rfloor} + (\gamma\rho - k^{-1/2} \|P_\phi \theta^*\|_2)^{2\lfloor s/2 \rfloor})/2 \\ &\simeq (\gamma\rho)^{\mathbf{1}\{s \text{ odd}\}} \cdot (\gamma|\rho| + k^{-1/2} \|P_\phi \theta^*\|_2)^{2\lfloor s/2 \rfloor}. \end{aligned} \quad (\text{F.30})$$

1407 Here, (i) holds by applying Lemma I.5 and \simeq denotes the equality that is up to a s -dependent
 1408 multiplicative constant.

1409 Putting together Eq. (F.29) and (F.30), we conclude that

$$\mathbb{E}_w[\langle w, \theta^* \rangle^s | \phi] \simeq (\gamma\rho)^{\mathbb{1}\{s \text{ odd}\}} (\gamma|\rho| + k^{-1/2} \|P_\phi \theta^*\|_2)^{s-1\{s \text{ odd}\}}. \quad (\text{F.31})$$

1410 In the sequel, we consider averaging over ϕ . From Eq. (F.31), we see that it suffices to consider
 1411 $\mathbb{E}_\phi[(\gamma|\rho| + k^{-1/2} \|P_\phi \theta^*\|_2)^r]$ for some $r \geq 2$. We alter the notation to facilitate some deferred
 1412 calculation. Consider $\mathbf{m} \subset [d]$ with constant size $|\mathbf{m}| = O(1)$ that does not scale with k or d . Now
 1413 define $\phi_{\mathbf{m}} \sim \text{Unif}\{\mathcal{S}_{k,\mathbf{m}}\}$, where $\mathcal{S}_{k,\mathbf{m}} = \{S \subset [d] : |S| = k, \mathbf{m} \subset S\}$. It is easily seen that this
 1414 definition covers previous definition of $\mathcal{S}_{k,m}$ by setting $\mathbf{m} = \{m\}$. We characterize the magnitude of
 1415 $\mathbb{E}_{\phi_{\mathbf{m}}}[\|P_{\phi_{\mathbf{m}}} \theta^*\|_2^r]$ from both sides as follows. For the lower bound, we have that

$$\begin{aligned} \mathbb{E}_{\phi_{\mathbf{m}}}[\|P_{\phi_{\mathbf{m}}} \theta^*\|_2^r] &= \mathbb{E}_{\phi_{\mathbf{m}}}\left[\left(\|\theta_{\mathbf{m}}^*\|_2^2 + \sum_{j \notin \mathbf{m}} |\theta_j^*|^2 \mathbb{1}\{j \in \phi_{\mathbf{m}}\}\right)^{r/2}\right] \\ &\geq \mathbb{E}_{\phi_{\mathbf{m}}}\left[\|\theta_{\mathbf{m}}^*\|_r^r + \sum_{j \notin \mathbf{m}} |\theta_j^*|^r \mathbb{1}\{j \in \phi_{\mathbf{m}}\}\right] \\ &\stackrel{(i)}{\gtrsim} (1 - k/d) \cdot \|\theta_{\mathbf{m}}^*\|_r^r + \frac{k}{d} \cdot \|\theta^*\|_r^r \\ &\stackrel{(ii)}{\gtrsim} \|\theta_{\mathbf{m}}^*\|_r^r + \frac{k}{d} \cdot k^{1-r/2} \cdot \|\theta^*\|_2^{r/2} \\ &= \|\theta_{\mathbf{m}}^*\|_r^r + \frac{k^2}{d} \cdot k^{-r/2}. \end{aligned}$$

1416 Here (i) holds by the fact that $\mathbb{E}[\mathbb{1}\{j \in \phi_{\mathbf{m}}\}] \simeq k/d$ for $j \notin \mathbf{m}$, and (ii) is a consequence of Jensen's
 1417 inequality. For the upper bound, we have that

$$\begin{aligned} \mathbb{E}_{\phi_{\mathbf{m}}}[\|P_{\phi_{\mathbf{m}}} \theta^*\|_2^r] &= \mathbb{E}_{\phi_{\mathbf{m}}}\left[\left(\|\theta_{\mathbf{m}}^*\|_2^2 + \sum_{j \notin \mathbf{m}} |\theta_j^*|^2 \cdot \mathbb{1}\{j \in \phi_{\mathbf{m}}\}\right)^{r/2}\right] \\ &\lesssim \mathbb{E}_{\phi_{\mathbf{m}}}\left[\|\theta_{\mathbf{m}}^*\|_r^r + \underbrace{\left(\sum_{j \notin \mathbf{m}} |\theta_j^*|^2 \cdot \mathbb{1}\{j \in \phi_{\mathbf{m}}\}\right)^{r/2}}_{|\phi_{\mathbf{m}} \cap \phi^* \setminus \{m\}| \text{ nonzero summands}}\right] \\ &\stackrel{\text{Jensen}}{\lesssim} \|\theta_{\mathbf{m}}^*\|_r^r + \mathbb{E}_{\phi_{\mathbf{m}}}\left[|\phi_{\mathbf{m}} \cap \phi^* \setminus \mathbf{m}|^{r/2-1} \cdot \left(\sum_{j \notin \mathbf{m}} |\theta_j^*|^r \cdot \mathbb{1}\{j \in \phi_{\mathbf{m}}\}\right)\right]. \quad (\text{F.32}) \end{aligned}$$

1418 Next, we apply Cauchy-Schwarz inequality as follows:

$$\begin{aligned} (\text{F.32}) &= \|\theta_{\mathbf{m}}^*\|_r^r + \mathbb{E}_{\phi_{\mathbf{m}}}\left[\sum_{j \notin \mathbf{m}} |\theta_j^*|^r \cdot \mathbb{1}\{j \in \phi_{\mathbf{m}}\}^2 \cdot |\phi_{\mathbf{m}} \cap \phi^* \setminus \mathbf{m}|^{r/2-1}\right] \\ &\leq \|\theta_{\mathbf{m}}^*\|_r^r + \mathbb{E}_{\phi_{\mathbf{m}}}\left[\sum_{j \notin \mathbf{m}} |\theta_j^*|^{2r} \cdot \mathbb{1}\{j \in \phi_{\mathbf{m}}\}\right]^{1/2} \cdot \mathbb{E}_{\phi_{\mathbf{m}}}\left[\sum_{j \notin \mathbf{m}} \mathbb{1}\{j \in \phi_{\mathbf{m}}\} \cdot |\phi_{\mathbf{m}} \cap \phi^* \setminus \mathbf{m}|^{r-2}\right]^{1/2} \\ &= \|\theta_{\mathbf{m}}^*\|_r^r + \left(\frac{k}{d} \cdot \sum_{j \notin \mathbf{m}} |\theta_j^*|^{2r}\right)^{1/2} \cdot \mathbb{E}_{\phi_{\mathbf{m}}}\left[|\phi_{\mathbf{m}} \cap \phi^* \setminus \mathbf{m}|^{r-1}\right]^{1/2} \\ &\stackrel{(i)}{\lesssim} \|\theta_{\mathbf{m}}^*\|_r^r + \left(\frac{k}{d} \cdot k^{1-r}\right)^{1/2} \cdot \left(\frac{k^2}{d}\right)^{(r-1)/2} \\ &= \|\theta_{\mathbf{m}}^*\|_r^r + \frac{k^2}{d} \cdot k^{-r/2}, \end{aligned}$$

1419 where (i) holds by $\mathcal{E}_{0,2r}$ and Lemma J.7. In conclusion, we have that $\mathbb{E}_{\phi_{\mathbf{m}}}[\|P_{\phi_{\mathbf{m}}} \theta^*\|_2^r] \simeq \|\theta_{\mathbf{m}}^*\|_r^r +$
 1420 $k^{-r/2} \cdot \delta$ given that $k = o(\sqrt{d})$. Combining this result with Eq. (F.31), we obtain that

$$\mathbb{E}_w[\langle w, \theta^* \rangle^s] \simeq (\gamma\rho)^{\mathbb{1}\{s \text{ odd}\}} (\gamma|\rho| + k^{-1/2} |\theta_{\mathbf{m}}^*| + k^{-1} \delta_{1/(s-1\{s \text{ odd}\})})^{s-1\{s \text{ odd}\}}$$

1421 where $\delta = k^2/d = o(1)$ and $\delta_r = \delta^{1/r}$. □

1422 **G Statistical Query Lower Bound for Sparse Signal Recovery**

1423 In this section, we provide a k^{s^*} sample complexity lower bound for the single index model with
 1424 k -sparse signal when querying a VSTAT oracle. The statistical query (SQ) framework was developed
 1425 in Feldman et al. (2017) and for completeness, we present essential definition and results here.

1426 **Definition G.1** (VSTAT Oracle). *Let D^* be the input distribution over domain \mathcal{X} . For a sample size*
 1427 *parameter $n > 0$, $\text{VSTAT}(D^*, n)$ oracle is the oracle that for any query function $h : \mathcal{X} \rightarrow [0, 1]$,*
 1428 *returns a value $v \in [p - \tau, p + \tau]$, where $p = \mathbb{E}_{x \sim D^*}[h(x)]$ and $\tau = \max\{t^{-1}, \sqrt{p(1-p)/n}\}$.*

1429 To define a key concept *statistical query dimension*, we first introduce the following notation.

1430 **Definition G.2** (Relative Pairwise Correlation). *Given two distributions $D_1, D_2 \in \Delta(\mathcal{X})$ and a*
 1431 *reference distribution $D \in \Delta(\mathcal{X})$,*

$$\chi_D(D_1, D_2) = \mathbb{E}_{x \sim D} \left[\frac{D_1(x)}{D(x)} \cdot \frac{D_2(x)}{D(x)} \right] - 1.$$

1432 **Definition G.3** (Statistical Dimension). *For $\bar{\gamma} > 0, \eta \in (0, 1)$, domain \mathcal{X} , a set of distributions \mathcal{D}*
 1433 *over \mathcal{X} , the **statistical dimension** $\text{SDA}(\mathcal{D}, \bar{\gamma}, \eta)$ of \mathcal{D} with average correlation $\bar{\gamma}$ and solution set*
 1434 *bound η is defined as the largest value m' such that there exists a reference distribution $D \in \Delta(\mathcal{X})$*
 1435 *and a finite set of distributions $\mathcal{D}_D \subseteq \mathcal{D}$ which can depend on the reference D with the following*
 1436 *property: for any solution $D^* \in \mathcal{D}$,*

1437 (i) $|\mathcal{D}_D \setminus \{D^*\}| \geq (1 - \eta)|\mathcal{D}_D|$;

1438 (ii) for any subset $\mathcal{D}'_D \subseteq \mathcal{D}_D \setminus \{D^*\}$ such that $|\mathcal{D}'_D| \geq |\mathcal{D}_D \setminus \{D^*\}|/m'$,

$$\frac{1}{|\mathcal{D}'_D|^2} \sum_{D_i, D_j \in \mathcal{D}'_D} \chi_D(D_i, D_j) \leq \bar{\gamma}.$$

1439 The above definition of the statistical dimension is a special case of the original Definition 3.1 in
 1440 Feldman et al. (2017) where we consider a search problem of *exact recovery* of the ground truth D^* .

Definition G.4 ((γ, β) -correlated Distributions). *We say that a set of m distributions $\mathcal{D} =$*
 $\{D_1, \dots, D_m\}$ over \mathcal{X} is (γ, β) -correlated relative to a reference distribution $D \in \Delta(\mathcal{X})$ if:

$$\chi_D(D_i, D_j) \leq \begin{cases} \beta & \text{for } i = j \in [m] \\ \gamma & \text{for } i \neq j \in [m]. \end{cases}$$

1441 The following lemma borrowed from Lemma 3.10 of Feldman et al. (2017) provides a lower bound on
 1442 the statistical dimension in terms of the (γ, β) -correlation property of the set of candidate distributions.

1443 **Lemma G.5.** *Given a set of candidate distributions \mathcal{D} that are (γ, β) -correlated with respect to a*
 1444 *reference distribution D , then for any $\gamma' > 0$ and $\eta > |\mathcal{D}|^{-1}$,*

$$\text{SDA}(\mathcal{D}, \gamma + \gamma', \eta) \geq \frac{(|\mathcal{D}| - 1)\gamma'}{\beta - \gamma}.$$

1445 The main result in the SQ framework is the following statement that relates the number of queries
 1446 required to the statistical dimension, which is borrowed from Theorem 3.2 of Feldman et al. (2017).

1447 **Lemma G.6.** *Let \mathcal{X} be a domain and \mathcal{D} be a set of candidate distributions over \mathcal{X} . For any $\bar{\gamma} > 0$*
 1448 *and $\eta \in (0, 1)$, Any randomized SQ algorithm that solves the problem of finding the input distribution*
 1449 *$D^* \in \mathcal{D}$ with probability at least $\alpha > \eta$ requires at least $(\alpha - \eta)/(1 - \eta) \cdot \text{SDA}(\mathcal{D}, \bar{\gamma}, \eta)$ calls to*
 1450 *the $\text{VSTAT}(D^*, (3\bar{\gamma})^{-1})$ oracle.*

1451 Our strategy for proving the lower bound is to first construct a set of candidate distributions \mathcal{D} that
 1452 are $(\omega(k^{-1}), \beta)$ -correlated with respect to reference distribution \mathbb{Q} with $\beta = D_{\chi^2}(\mathbb{P}_{\theta^*} \parallel \mathbb{Q})$ and
 1453 $|\mathcal{D}|$ exponentially large. Then by Lemma G.5 and Lemma G.6, we can derive the desired hardness
 1454 result. It remains to construct the set of candidate distributions \mathcal{D} that are $(\omega(k^{-1}), \beta)$ -correlated
 1455 with respect to \mathbb{Q} . To this end, we introduce the following result on the packing number of k -sparse
 1456 vectors.

1457 **Lemma G.7** (Packing Number for k -Sparse Vectors). Define $\rho(u, v) = |\langle u, v \rangle|$. Let packing number
 1458 $\mathcal{M}_\rho(d, k, t)$ be the maximal cardinality of the set of k -sparse vectors in \mathbb{S}^{d-1} such that $\rho(u, v) < t$
 1459 for any $u \neq v$ in the set. We have for any $t \in (1/k, 1)$ that

$$\mathcal{M}_\rho(d, k, t) \geq \frac{1}{2} \cdot \exp\left(\frac{\min\{(d-k)t^2, 3kt\}}{8}\right).$$

1460 With all these ingredients in place, we are ready to prove the main theorem.

1461 *Proof of Theorem 5.4.* Let us pick parameter $\kappa_d \in ((\log d)^2, k/4)$ that scales with d and set

$$t \geq \max\left\{\sqrt{\frac{\kappa_d}{d-k}}, \frac{\kappa_d}{3k}\right\} \in (1/k, 1/2). \quad (\text{G.1})$$

1462 Note that $t \in (1/k, 1/2)$ is able to hold by our choice of κ_d and condition that $\omega((\log d)^2) \leq k \leq d/2$.
 1463 In this vein, we can pick \mathcal{D} to be the maximal set of distributions \mathbb{P}_θ for some k -sparse vectors
 1464 $\theta \in \mathbb{S}^{d-1}$ satisfying $\rho(\theta, \theta') < t$ for any $\theta \neq \theta'$ in the set. It follows from plugging (G.1) into
 1465 Lemma G.7 that $|\mathcal{D}| \geq \exp(\kappa_d/8)/2$, which is super polynomially large in d for our choice of κ_d .

1466 Next, we configure the remaining parameters in Lemma G.5 and Lemma G.6. We choose the reference
 1467 distribution to be \mathbb{Q} , in which the covariate z is independent of the output y . For β , we note that

$$\chi_{\mathbb{Q}}(\mathbb{P}_\theta, \mathbb{P}_\theta) = D_{\chi^2}(\mathbb{P}_\theta \| \mathbb{Q}) = O(1),$$

1468 which is a constant independent of θ due to the rotational invariance of the likelihood ratio with
 1469 respect to θ . Thus, we define this quantity as B can just set $\beta = D_{\chi^2}(\mathbb{P}_{\theta^*} \| \mathbb{Q}) = B$. For γ , we note
 1470 that for any two $\mathbb{P}_\theta, \mathbb{P}_{\theta'}$ in \mathcal{D} for $\theta \neq \theta'$,

$$\begin{aligned} |\chi_{\mathbb{Q}}(\mathbb{P}_\theta, \mathbb{P}_{\theta'})| &= \left| \mathbb{E}_{x \sim \mathbb{Q}} \left[\frac{\mathbb{P}_\theta(x)}{\mathbb{Q}(x)} \cdot \frac{\mathbb{P}_{\theta'}(x)}{\mathbb{Q}(x)} \right] - 1 \right| \\ &= \left| \mathbb{E}_{x \sim \mathbb{Q}} \left[\left(1 + \sum_{s \geq s^*} \zeta_s(y) h_s(\langle \theta, z \rangle) \right) \cdot \left(1 + \sum_{s' \geq s^*} \zeta_{s'}(y) h_{s'}(\langle \theta', z \rangle) \right) \right] - 1 \right| \\ &= \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \cdot |\langle \theta, \theta' \rangle|^s \leq \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \cdot t^s \leq \mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y)^2] \cdot t^{s^*} + \frac{t^{s^*+1}}{1-t}. \end{aligned}$$

1471 Here, the third equality follows from the fact that only when $s = s'$, the cross term
 1472 $\mathbb{E}_{\mathbb{Q}}[h_s(\langle \theta, z \rangle) h_{s'}(\langle \theta', z \rangle)]$ is non-zero. In particular, by the property of the Gaussian noise opera-
 1473 tor introduced in (B.3), we have that $\mathbb{E}_{\mathbb{Q}}[h_s(\langle \theta, z \rangle) h_{s'}(\langle \theta', z \rangle)] = \langle \theta, \theta' \rangle^s < t^s$. For the last inequality
 1474 above, we simply use the fact that $\mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \leq 1$ for any s (Damian et al., 2024) and $t < 1$. Now,
 1475 we conclude that

$$|\chi_{\mathbb{Q}}(\mathbb{P}_\theta, \mathbb{P}_{\theta'})| \leq \left(\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y)^2] + \frac{t}{1-t} \right) \cdot t^{s^*} \leq (\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y)^2] + 1) \cdot t^{s^*}.$$

1476 We thus set $\gamma' = \gamma = (\mathbb{E}_{\mathbb{Q}}[\zeta_{s^*}(y)^2] + 1) \cdot t^{s^*} = \Theta(t^{s^*})$. Finally, we set $\eta = 1/3$ and $\alpha = 2/3$. Then
 1477 all the conditions in both Lemma G.5 and Lemma G.6 are satisfied and we have

$$\text{SDA}(\mathcal{D}, 2\gamma, 1/3) \geq \frac{(|\mathcal{D}| - 1)\gamma}{\beta - \gamma} \geq \frac{|\mathcal{D}|\gamma}{2\beta} \geq \frac{\gamma \exp(\kappa_d/8)}{4\beta}.$$

1478 Lastly, recall that we have $|\langle \theta, \theta' \rangle| \leq t$ for any $\theta \neq \theta'$ in \mathcal{D} , which means that in order to achieve
 1479 alignment at least $2t$ with the true signal θ^* , we need to *exactly* identify the distribution \mathbb{P}_{θ^*} . Con-
 1480 sequently, by Lemma G.6, we have that any randomized SQ algorithm that solves the problem of
 1481 achieving alignment $2t$ with probability at least $2/3$ requires at least $\gamma \exp(\kappa_d/8)/(8B)$ calls to the
 1482 $\text{VSTAT}(\mathbb{P}_{\theta^*}, (6\gamma)^{-1})$ oracle.

1483 **Simplification of the lower bound.** To simplify the lower bound, let us take $\kappa_d = (\log d)^c/2$ for
 1484 some constant $c > 2$. Thus, the alignment $2t$ is upper bounded by

$$2t \leq \begin{cases} \tilde{\omega}(k^{-1}) & \text{if } k < \sqrt{d} \\ \tilde{\omega}(d^{-1/2}) & \text{if } k \geq \sqrt{d} \end{cases},$$

1485 where $\tilde{\omega}(\cdot)$ hides some poly-logarithmic factors. The number of queries is still super polynomially
 1486 large in d . Following from (G.1), we can safely set

$$t = \begin{cases} (\log d)^c/k & \text{if } (\log d)^2 < k < \sqrt{d(\log d)^c} \\ \sqrt{(\log d)^c/d} & \text{if } \sqrt{d(\log d)^c} \leq k \leq d/2 \end{cases},$$

1487 Hence, the number of sample

$$(6\gamma)^{-1} = \frac{t^{-s^*}}{6} \simeq \begin{cases} \frac{k^{s^*}}{(\log d)^{cs^*}} & \text{if } (\log d)^2 < k < \sqrt{d(\log d)^c} \\ \frac{d^{s^*/2}}{6(\log d)^{cs^*/2}} & \text{if } \sqrt{d(\log d)^c} \leq k \leq d/2 \end{cases}.$$

1488 Hence, we have established the desired lower bound on the sample complexity. \square

1489 *Proof of Lemma G.7.* We use the probability method to prove the existence of a set of k -sparse vectors
 1490 in \mathbb{S}^{d-1} with the desired property. We i.i.d. sample m vectors $\omega^{(1)}, \dots, \omega^{(m)}$ from the following
 1491 distribution:

$$\omega : \phi \sim \text{Unif}(\mathcal{S}_k), \quad \omega_j = \begin{cases} \frac{1}{\sqrt{k}}, & \text{w.p. } \frac{1}{2} \text{ if } j \in \phi \\ -\frac{1}{\sqrt{k}}, & \text{w.p. } \frac{1}{2} \text{ if } j \in \phi, \\ 0, & j \notin \phi. \end{cases} \quad j \in [d].$$

1492 where we recall that \mathcal{S}_k is the set of all size- k subsets in $[d]$. Since each $\omega^{(i)}$ is i.i.d. sampled, we can
 1493 equivalently view $\langle \omega^{(i)}, \omega^{(j)} \rangle$ for $i \neq j$ as a random variable sampled from the following distribution:

$$\langle \omega^{(i)}, \omega^{(j)} \rangle \stackrel{d}{=} \frac{R_X}{k}, \quad \text{where } R_X = r_1 + \dots + r_X, \quad X \sim \text{Hypergeometric}(d, k, k), \quad (\text{G.2})$$

1494 where r_1, r_2, \dots are i.i.d. Rademacher random variables. Let us consider random variable W dis-
 1495 tributed as

$$W \stackrel{d}{=} \frac{R_Y}{k}, \quad \text{where } R_Y = r_1 + \dots + r_Y, \quad Y \sim \text{Binomial}\left(k, \frac{k}{d-k}\right). \quad (\text{G.3})$$

1496 We will invoke the following fact on the tail probability regarding the above two random variables.

1497 **Proposition G.8.** For R_X and R_Y defined in (G.2) and (G.3), respectively, we have that $\mathbb{P}(R_X \geq$
 1498 $t) \leq 2\mathbb{P}(R_Y \geq t)$ for any $t > 1$.

1499 The proof of the proposition is deferred to the end of the proof. Thus, it suffices to study the tail
 1500 probability of W . Note that $W \stackrel{d}{=} \sum_{j=1}^k w_j$ where w_j are i.i.d. sampled from

$$w_j = \begin{cases} \frac{1}{k}, & \text{w.p. } \frac{k}{2(d-k)} \\ -\frac{1}{k}, & \text{w.p. } \frac{k}{2(d-k)} \\ 0, & \text{w.p. } 1 - \frac{k}{d-k} \end{cases}, \quad j \in [k].$$

1501 where $\mathbb{E}[w_j] = 0$ and $\mathbb{E}[w_j^2] = (k(d-k))^{-1}$. Hence, we can apply the Bernstein inequality to obtain
 1502 that for any $t > 1/k$,

$$\begin{aligned} \mathbb{P}(\langle \omega^{(i)}, \omega^{(j)} \rangle \geq t) &\leq 2\mathbb{P}(W \geq t) \leq 2 \exp\left(-\frac{k(t/k)^2/2}{(k(d-k))^{-1} + t/(3k^2)}\right) \\ &= 2 \exp\left(-\frac{k^2 t^2}{2k^2/(d-k) + 2kt/3}\right) \leq 2 \exp\left(-\min\left\{\frac{(d-k)t^2}{4}, \frac{3kt}{4}\right\}\right). \end{aligned}$$

1503 Suppose we randomly sample m i.i.d. $\omega^{(i)}$ from the same distribution. Then the probability that all
 1504 such pair $|\langle \omega^{(i)}, \omega^{(j)} \rangle| < t$ for $t > 1/k$ is lower bounded by

$$\begin{aligned} \mathbb{P}(|\langle \omega^{(i)}, \omega^{(j)} \rangle| < t, \forall i \neq j) &\geq 1 - m^2 \cdot 2\mathbb{P}(\langle \omega^{(i)}, \omega^{(j)} \rangle \geq t) \\ &\geq 1 - 4m^2 \cdot \exp\left(-\min\left\{\frac{(d-k)t^2}{4}, \frac{3kt}{4}\right\}\right). \end{aligned}$$

1505 Ensuring that the probability is nonzero will give us a valid construction of the set \mathcal{D} . Therefore, there
 1506 must exist a \mathcal{D} satisfying $|\langle \omega^{(i)}, \omega^{(j)} \rangle| < t$ for any $i \neq j$ and with size

$$|\mathcal{D}| \geq \frac{1}{2} \cdot \exp\left(\frac{\min\{(d-k)t^2, 3kt\}}{8}\right).$$

1507 Hence, we complete the proof. \square

1508 Next, we aim to present the proof of Proposition G.8. To proceed, let us introduce the definition of
 1509 stochastic dominance.

1510 **Definition G.9** (Stochastic Dominance). *For any real-valued random variable X and Y , we say that*
 1511 *X is stochastically dominated by Y , denoted by $X \stackrel{\text{s.t.}}{\leq} Y$, if $\mathbb{P}(X \geq t) \leq \mathbb{P}(Y \geq t)$ for every t .*

1512 The following result is from Theorem A, Chapter 2 of Szekli (2012).

1513 **Proposition G.10.** *We have $X \stackrel{\text{s.t.}}{\leq} Y$ if and only if there exists a coupling (\hat{X}, \hat{Y}) with $\text{law}(\hat{X}) =$
 1514 $\text{law}(X)$ and $\text{law}(\hat{Y}) = \text{law}(Y)$ such that $\hat{X} \leq \hat{Y}$ almost surely.*

1515 **Proposition G.11** (Theorem 1.1, Klenke and Mattner (2010)). *Hypergeometric(d, k, k) $\stackrel{\text{s.t.}}{\leq}$*
 1516 *Binomial($k, k/(d-k)$).*

1517 Another way to think of the problem is that Hypergeometric(d, k, k) corresponds to the number of
 1518 times a black ball is drawn when sampling for k times from an urn with $d-k$ white ball and k black
 1519 ball without replacement, while Binomial($k, k/(d-k)$) corresponds to sampling in the same urn but
 1520 with replacement. We claim the following fact on the tail probability of sum of Rademacher random
 1521 variables.

1522 **Proposition G.12** (Sum of Rademacher Random Variables). *Let r_1, r_2, \dots be i.i.d. Rademacher*
 1523 *random variables. Let $R_l = r_1 + \dots + r_l$ for $l = 1, 2, \dots$. Let $p_l(\cdot)$ be the probability mass function*
 1524 *of B_l . Then the following holds for any $l = 1, 2, \dots$:*

- 1525 1. p_l is symmetric and supported on the set of odd integers if l is odd, and supported on the set
 1526 of even integers if l is even.
- 1527 2. For $i \in \text{supp}(p_l)$ and $i \geq 0$, $p_l(i)$ is a non-increasing function of i .
- 1528 3. $\mathbb{P}(R_l \geq t) \leq \mathbb{P}(R_{l+2} \geq t)$ for any $t > 1$.
- 1529 4. $\mathbb{P}(R_l \geq t) \leq 2\mathbb{P}(R_{l+1} \geq t)$ for any $t > 1$.
- 1530 5. $\mathbb{P}(R_l \geq t) \leq 2\mathbb{P}(R_{l+l'} \geq t)$ for any $l \geq 1$ and $l' \geq 1$.

1531 *Proof of Proposition G.12.* The first claim is immediate from the symmetry of the Rademacher
 1532 random variables and the fact that the sum of an odd number of Rademacher random variables is odd,
 1533 while the sum of an even number of Rademacher random variables is even. For the second claim, we
 1534 note that

$$p_l(i) = 2^{-l} \cdot \binom{l}{(i+l)/2}, \quad i \in \text{supp}(p_l),$$

1535 which is a non-increasing function for $i \geq 0$. For the third claim, we let $t^* = 2\lceil t/2 \rceil$ if l is even and
 1536 $t^* = 2\lceil (t-1)/2 \rceil + 1$ if l is odd. In other words, $t^* = \min\{\tau \in \text{supp}(p_l) : \tau \geq t\}$. Then we have
 1537 that

$$\begin{aligned} \mathbb{P}(R_{l+2} \geq t) &= \mathbb{P}(R_l \geq t^* + 2) + \mathbb{P}(R_l = t^*) \cdot \mathbb{P}(r_{l+1} + r_{l+2} \geq 0) \\ &\quad + \mathbb{P}(R_l = t^* - 2) \cdot \mathbb{P}(r_{l+1} + r_{l+2} = 2) \\ &= \mathbb{P}(R_l \geq t^*) + (\mathbb{P}(R_l = t^* - 2) - \mathbb{P}(R_l = t^*)) \cdot \mathbb{P}(r_{l+1} + r_{l+2} = 2) \\ &\geq \mathbb{P}(R_l \geq t^*) = \mathbb{P}(R_l \geq t). \end{aligned}$$

1538 where in the first equality we use the fact that $r_{l+1} + r_{l+2}$ is supported on $\{-2, 0, 2\}$ and in the
 1539 second equality we use the symmetric property of the distribution of $r_{l+1} + r_{l+2}$. The last inequality

1540 follows from the monotonicity of the probability mass function of R_l for $t^* - 2 \geq 0$ when $t > 1$. For
 1541 the forth claim, we similarly have that

$$\mathbb{P}(R_{l+1} \geq t) \geq \mathbb{P}(R_l \geq t^*) - \mathbb{P}(R_l = t^*) \cdot \mathbb{P}(r_{l+1} = -1) \geq \frac{1}{2} \mathbb{P}(R_l \geq t^*) = \frac{1}{2} \mathbb{P}(R_l \geq t).$$

1542 The last claim follows from a combination of the third and forth claims where

$$\mathbb{P}(R_{l+\nu} \geq t) \geq \frac{1}{2} \mathbb{P}(R_{l+2\lfloor \nu/2 \rfloor} \geq t) \geq \frac{1}{2} \mathbb{P}(R_{l+2\lfloor \nu/2 \rfloor - 2} \geq t) \geq \dots \geq \frac{1}{2} \mathbb{P}(R_l \geq t).$$

1543 Hence, the proof is complete. \square

1544 Next, we proceed to the proof of Proposition G.8.

1545 *Proof of Proposition G.8.* By Proposition G.11 and Proposition G.10, there exists a coupling \hat{X}, \hat{Y}
 1546 with $\text{law}(\hat{X}) = \text{law}(X)$ and $\text{law}(\hat{Y}) = \text{law}(Y)$ such that $\hat{X} \leq \hat{Y}$ almost surely where $X \sim$
 1547 Hypergeometric(d, k, k) and $Y \sim \text{Binomial}(k, k/(d-k))$.

1548 Consider i.i.d. Rademacher random variables r_1, r_2, \dots, r_k . Let $R_l = r_1 + \dots + r_l$ for $l = 1, 2, \dots$
 1549 Since $R_{\hat{X}} = r_1 + \dots + r_{\hat{X}} \mid \hat{X} \stackrel{d}{=} 2\text{Binomial}(\hat{X}, 1/2) - L$ and $R_{\hat{Y}} = r_1 + \dots + r_{\hat{Y}} \mid \hat{Y} \stackrel{d}{=} 2\text{Binomial}(\hat{Y}, 1/2) - L$
 1550 for the coupling (\hat{X}, \hat{Y}) with $\hat{X} \leq \hat{Y}$, we consider the conditional random
 1551 variable

$$r_{\hat{X}+i} \mid (R_{\hat{X}+i-1}, \hat{X}, \hat{Y}) = r_{\hat{X}+i} = \begin{cases} 1, & \text{w.p. } 1/2 \\ -1, & \text{w.p. } 1/2 \end{cases}, \quad i = 1, 2, \dots, \hat{Y} - \hat{X}$$

1552 The equality holds by the i.i.d. property of these Rademacher random variables. From the distributional
 1553 perspective, the distribution of $R_{\hat{Y}}$ is obtained by conducting convolution with the Rademacher
 1554 distribution for $\hat{Y} - \hat{X}$ times on the distribution of $R_{\hat{X}}$. Invoking Proposition G.12, we directly
 1555 conclude that $\mathbb{P}(R_{\hat{Y}} \geq t \mid \hat{X}, \hat{Y}) \geq \mathbb{P}(R_{\hat{X}} \geq t \mid \hat{X}, \hat{Y})/2$ for any $t > 1$ and $\hat{Y} \geq \hat{X}$. As $\hat{Y} \geq \hat{X}$
 1556 holds almost surely, by the law of total probability, we arrive at the conclusion that $\mathbb{P}(R_{\hat{Y}} \geq t) \geq$
 1557 $\mathbb{P}(R_{\hat{X}} \geq t)/2$ for any $t > 1$. \square

1558 H Supporting Lemmas on Moment Calculations

1559 **Lemma H.1** (First moment). *Under Assumption 4.1, for any $s \geq 0$, it holds for any $y \in \mathbb{R}$ and*
 1560 *$w, \theta \in \mathbb{S}^{d-1}$ that*

$$\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, \langle w, z \rangle) z \cdot h_s(\langle \theta, z \rangle)] = \sqrt{s+1} \cdot \hat{\psi}_{s+1}(y) \cdot \langle w, \theta \rangle^s w + \sqrt{s} \cdot \hat{\psi}_{s-1}(y) \cdot \langle w, \theta \rangle^{s-1} \theta,$$

1561 *in the L^2 sense over the marginal distribution of y under \mathcal{Q} .*

1562 *Proof of Lemma H.1.* For convenience, we denote $\rho := \langle w, \theta^* \rangle$. We claim the following identities:

$$\begin{aligned} & \mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, w^\top z) z \cdot h_s(\theta^{*\top} z)] \\ &= \mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, w^\top z) \cdot \theta^{*\top} z \cdot h_s(\theta^{*\top} z)] \cdot \theta^* + \mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, w^\top z) \cdot h_s(\theta^{*\top} z) \cdot P_{\theta^*}^\perp z] \quad (\text{H.1}) \\ &= \underbrace{\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, w^\top z) \cdot \theta^{*\top} z \cdot h_s(\theta^{*\top} z)]}_{(\text{I})} \cdot \frac{\theta^* - \rho w}{1 - \rho^2} + \underbrace{\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, w^\top z) \cdot w^\top z \cdot h_s(\theta^{*\top} z)]}_{(\text{II})} \cdot \frac{w - \rho \theta^*}{1 - \rho^2}. \end{aligned}$$

1563 Here, in the first identity, we project z in the direction of θ^* and the orthogonal complement of θ^* ,
 1564 where $P_{\theta^*}^\perp = I - \theta^* \theta^{*\top}$ is the projection operator onto the orthogonal complement of θ^* . To see how
 1565 the second identity holds, we first look at the second term $\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, w^\top z) \cdot h_s(\theta^{*\top} z) \cdot P_{\theta^*}^\perp z]$.
 1566 For each direction v orthogonal to both θ^* and w , we have

$$\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(y, w^\top z) \cdot h_s(\theta^{*\top} z) \cdot \langle P_{\theta^*}^\perp z, v \rangle] = \mathbb{E}_{z \sim \mathcal{N}_d, x \sim \mathcal{N}} [\psi(y, w^\top z) \cdot h_s(\theta^{*\top} z) \cdot x] = 0.$$

1567 Also, by projection $P_{\theta^*}^\perp z$ is always orthogonal to θ^* . Thus, the only direction left for consideration is
 1568 $v = (w - \rho\theta^*)/\sqrt{1 - \rho^2}$, for which we have

$$\begin{aligned} & \mathbb{E}_{z \sim \mathcal{N}_d} \left[\psi(y, w^\top z) \cdot h_s(\theta^{*\top} z) \cdot \langle P_{\theta^*}^\perp z, v \rangle \right] \cdot v \\ &= \mathbb{E}_{z \sim \mathcal{N}_d} \left[\psi(y, w^\top z) \cdot h_s(\theta^{*\top} z) \cdot \frac{w^\top z - \rho\theta^{*\top} z}{\sqrt{1 - \rho^2}} \right] \cdot \frac{w - \rho\theta^*}{\sqrt{1 - \rho^2}} \\ &= \mathbb{E}_{z \sim \mathcal{N}_d} \left[\psi(y, w^\top z) \cdot h_s(\theta^{*\top} z) \cdot (w^\top z - \rho\theta^{*\top} z) \right] \cdot \frac{w - \rho\theta^*}{1 - \rho^2}. \end{aligned} \quad (\text{H.2})$$

1569 Plugging Eq. (H.2) into the second term of line 2 in Eq. (H.1), we thus have the last identity in
 1570 Eq. (H.1). Next, we analyze terms (I) and (II) in Eq. (H.1). For our convenience, we define U_ρ as the
 1571 Gaussian noise operator such that

$$U_\rho \psi(y, x) = \mathbb{E}_{x' \sim \mathcal{N}} \left[\psi(y, \rho x + \sqrt{1 - \rho^2} x') \right].$$

1572 For term (I), we have by the definition of U_ρ that

$$\begin{aligned} (\text{I}) &= \mathbb{E}_{x \sim \mathcal{N}} [U_\rho \psi(y, x) \cdot x \cdot h_s(x)] \\ &= \mathbb{E}_{z \sim \mathcal{N}_d} \left[\sqrt{s+1} \cdot U_\rho \psi(y, x) \cdot h_{s+1}(x) + \sqrt{s} \cdot U_\rho \psi(y, x) \cdot h_{s-1}(x) \right] \\ &\stackrel{L^2(\mathbb{Q})}{=} \sqrt{s+1} \cdot \hat{\psi}_{s+1}(y) \cdot \rho^{s+1} + \sqrt{s} \cdot \hat{\psi}_{s-1}(y) \cdot \rho^{s-1}. \end{aligned} \quad (\text{H.3})$$

1573 where the second line follows from the recurrence relation of the Hermite polynomials in Eq. (B.1),
 1574 and the last line follows from the property of the Gaussian noise operator in Eq. (B.3). Similarly for
 1575 term (II), we have

$$\begin{aligned} (\text{II}) &= \mathbb{E}_{x \sim \mathcal{N}} [U_\rho (\psi(y, x)x) \cdot h_s(x)] \\ &\stackrel{L^2(\mathbb{Q})}{=} \rho^s \cdot \mathbb{E}_{x \sim \mathcal{N}} [\psi(y, x) \cdot x \cdot h_s(x)] \\ &\stackrel{L^2(\mathbb{Q})}{=} \rho^s \cdot \left(\sqrt{s+1} \cdot \hat{\psi}_{s+1}(y) + \sqrt{s} \cdot \hat{\psi}_{s-1}(y) \right), \end{aligned} \quad (\text{H.4})$$

1576 where in the last line we borrow the calculation in Eq. (H.3) by letting $\rho = 1$. Plugging Eq. (H.3)
 1577 and (H.4) into Eq. (H.1), we hence have

$$\begin{aligned} (\text{H.1}) &\stackrel{L^2(\mathbb{Q})}{=} \left(\sqrt{s+1} \cdot \hat{\psi}_{s+1}(y) \cdot \rho^{s+1} + \sqrt{s} \cdot \hat{\psi}_{s-1}(y) \cdot \rho^{s-1} \right) \cdot \frac{\theta^* - \rho w}{1 - \rho^2} \\ &\quad + \rho^s \cdot \left(\sqrt{s+1} \cdot \hat{\psi}_{s+1}(y) + \sqrt{s} \cdot \hat{\psi}_{s-1}(y) \right) \cdot \frac{w - \rho\theta^*}{1 - \rho^2} \\ &= \sqrt{s+1} \cdot \hat{\psi}_{s+1}(y) \cdot \rho^s w + \sqrt{s} \cdot \hat{\psi}_{s-1}(y) \cdot \rho^{s-1} \theta^*, \end{aligned}$$

1578 which completes the proof. \square

1579 An implication of the previous lemma is that

$$\mathbb{E}_{\mathbb{Q}} [h_{s^*}(\langle \theta^*, z \rangle) \cdot \sigma'(\langle z, \theta \rangle) \cdot \langle z, \theta^* \rangle] = s \cdot \hat{\sigma}^{(s^*)} \cdot \langle \theta^*, \theta \rangle^{s^*-1} + \sqrt{(s+1)(s+2)} \cdot \hat{\sigma}^{(s^*+2)} \cdot \langle \theta^*, \theta \rangle^{s^*+1},$$

1580 where we take $\hat{\sigma}^{(s)}$ as the s -th normalized Hermite coefficient of σ . Here, we take $\psi(y, x)$ as $\sigma'(x)$
 1581 and thus $\hat{\psi}_s(y) = \sqrt{s+1} \cdot \hat{\sigma}^{(s+1)}$.

1582 **Lemma H.2** (Decomposition of first order moment). *Suppose that ψ follows Assumption 4.1 and*

$$g = \frac{1}{nL} \sum_{i=1}^n \sum_{l=1}^L (\psi(y_i, \langle w_l, z_i \rangle) \cdot z_i - \hat{\psi}_1(y_i) \cdot w_l),$$

1583 where $(z_i, y_i) \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{\theta^*}$ and $\{w_l\}_{l \leq L}$ is fixed. Then it holds that

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{\theta^*}} [g] &= \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \cdot \hat{\psi}_{s-1}(y)] \cdot \frac{\sqrt{s}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^{s-1} \cdot \theta^* \\ &\quad + \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \cdot \hat{\psi}_{s+1}(y)] \cdot \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^s \cdot w_l. \end{aligned}$$

1584 *Proof of Lemma H.2.* Applying a change of measure from \mathbb{P}_{θ^*} to \mathbb{Q} and invoking Eq. (2.2), we get

$$\begin{aligned}\mathbb{E}_{\mathbb{P}_{\theta^*}}[g] &= \frac{1}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{P}_{\theta^*}} \left[\psi(y, \langle w_l, z \rangle) \cdot z - \widehat{\psi}_1(y) \cdot w_l \right] \\ &= \frac{1}{L} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}} \left[\psi(y, \langle w_l, z \rangle) z \cdot \left(1 + \sum_{s \geq s^*} \zeta_s(y) h_s(\langle \theta^*, z \rangle) \right) - \widehat{\psi}_1(y) w_l \right],\end{aligned}\quad (\text{H.5})$$

1585 Note that for $s = 0$, we have for the first term in the summation of Eq. (H.5) that

$$L^{-1} \sum_{l=1}^L \mathbb{E}_{\mathbb{Q}}[\psi(y, \langle w_l, z \rangle) z] = \mathbb{E}_{\mathbb{Q}}[\widehat{\psi}_1(y)] \cdot \frac{1}{L} \sum_{l=1}^L w_l,$$

1586 which is cancelled out by the debiasing term in the algorithm. Applying the result of Lemma H.1 to
1587 the remaining terms in Eq. (H.5) yields

$$\begin{aligned}\mathbb{E}_{\mathbb{P}_{\theta^*}}[g] &= \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \widehat{\psi}_{s+1}(y)] \cdot \frac{\sqrt{s+1}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^s \cdot w_l \\ &\quad + \sum_{s \geq s^*} \mathbb{E}_{\mathbb{Q}}[\zeta_s(y) \cdot \widehat{\psi}_{s-1}(y)] \cdot \frac{\sqrt{s}}{L} \sum_{l=1}^L \langle w_l, \theta^* \rangle^{s-1} \cdot \theta^*,\end{aligned}\quad (\text{H.6})$$

1588 where the $\mathbb{E}_{\mathbb{Q}}[\widehat{\psi}_1(y)] \cdot \frac{1}{L} \sum_{l=1}^L w_l$ term from Lemma H.1 with $s = 0$ is cancelled out by the debiasing
1589 term in the algorithm. \square

1590 **Lemma H.3** (Second moment on nice event). *Suppose $\psi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ satisfies the quadruple-*
1591 *integrable and high-pass assumptions in Assumption 4.1. Let s^* be the generative exponent defined*
1592 *in Definition 2.1. Suppose $\mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \leq C$ for some universal $C = O(1)$ and for all $s \geq s^*$. For any*
1593 *$w, w', \theta^*, v \in \mathbb{S}^{d-1}$ where either $v = \theta^*$ or $\langle v, \theta^* \rangle = 0$ in the non-sparse case, and either $v = e_j$ for*
1594 *$j \in \text{supp}(\theta^*)$ or $v = e_j$ for $j \notin \text{supp}(\theta^*)$ in the sparse case, suppose that*

$$\max\{|\langle v, w \rangle|, |\langle v, w' \rangle|\} \leq \epsilon_0, \quad \max\{|\langle \theta^*, w \rangle|, |\langle \theta^*, w' \rangle|\} \leq \epsilon, \quad |\langle w, w' \rangle| \leq \epsilon_1$$

1595 for some $\epsilon, \epsilon_0, \epsilon_1$ such that $4\epsilon s^* \epsilon < 1/2$. Then, we have for $s^* \geq 2$ that

$$\begin{aligned}\mathbb{E}_{\mathbb{Q}} \left[\psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \langle v, z \rangle^2 \cdot \left(1 + \sum_{s=s^*}^{\infty} \zeta_s(y) h_s(\langle \theta^*, z \rangle) \right) \right] \\ \lesssim \epsilon_1^{s^*-1} \cdot \left(1 + \frac{\epsilon^2}{\epsilon_1} + \left(\frac{\epsilon^2}{\epsilon_1} \right)^{s^*-1} \cdot \epsilon + \mathbf{1}(v \perp \theta^*) \cdot \left(\frac{\epsilon^2}{\epsilon_1} \right)^{s^*-2} \cdot \frac{\epsilon_0^2}{\epsilon_1} \cdot (\epsilon^2 + \epsilon \cdot \mathbf{1}(s^* \geq 4)) \right),\end{aligned}$$

1596 and for $s^* = 1$, the bound is $O(1)$. Here, \lesssim hides constants that only depend on s^* , $\mathbb{E}_{\mathbb{Q}}[\psi(x, y)^4]$
1597 and C .

1598 *Proof.* Using the results from Proposition I.1, we have that

$$\begin{aligned}h_s(\langle \theta^*, z \rangle) \langle v, z \rangle^2 &= \sqrt{(s+2)(s+1)} \cdot \mathbf{h}_{s+2}(z) [(\theta^*)^{\otimes s} \otimes v^{\otimes 2}] + \mathbf{h}_s(z) [(\theta^*)^{\otimes s}] \\ &\quad + 2s \cdot \mathbf{h}_s(z) [(\theta^*)^{\otimes s-1} \otimes v]^\top \cdot \langle \theta^*, v \rangle + \sqrt{s(s-1)} \cdot \mathbf{h}_{s-2}(z) [(\theta^*)^{\otimes s-2}] \cdot \langle \theta^*, v \rangle^2.\end{aligned}\quad (\text{H.7})$$

1599 Thus, we only need to focus on these degree terms in $\psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle)$. Our goal is to
1600 compute the following quantity, which we denoted by F :

$$\begin{aligned}F &= \mathbb{E}_{\mathbb{Q}} \left[\psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \langle v, z \rangle^2 \cdot \left(1 + \sum_{s=s^*}^{\infty} \zeta_s(y) h_s(\langle \theta^*, z \rangle) \right) \right] \\ &= \mathbb{E}_{\mathbb{Q}} \left[\psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \langle v, z \rangle^2 \right] \\ &\quad + \sum_{s=s^*}^{\infty} \left| \mathbb{E}_{\mathbb{Q}} \left[\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \langle v, z \rangle^2 h_s(\langle \theta^*, z \rangle) \right] \right|.\end{aligned}\quad (\text{H.8})$$

1601 Here, for the term corresponding to $s = 0$ in Eq. (H.8), we plug in Eq. (H.7) and have by Lemma I.3
 1602 that

$$\begin{aligned} & \left| \mathbb{E}_{\mathbb{Q}} [\psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \langle v, z \rangle^2] \right| \\ &= \left| \mathbb{E}_{\mathbb{Q}} [\psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) (\sqrt{2} \cdot \mathbf{h}_2(z) [v^{\otimes 2}] + 1)] \right| \\ &\lesssim \epsilon_1^{(s^*-2) \vee 0} \cdot \epsilon_0^{2 \wedge c_0} \cdot \epsilon^{(2-c_0) \vee 0} + \epsilon_1^{s^*-1} \lesssim \epsilon_1^{(s^*-2) \vee 0} \cdot \epsilon_0^2 + \epsilon_1^{s^*-1} \lesssim \mathbb{1}(s^* = 1) + \epsilon_1^{s^*-2} \epsilon_0^2 + \epsilon_1^{s^*-1}. \end{aligned}$$

1603 Here, to use Lemma I.3, for $\mathbf{h}_2(z) [v^{\otimes 2}]$ we take test tensor $T_2 = v^{\otimes 2}$ and set $c_0 = 2$. The last line
 1604 also holds by using the Cauchy-Schwarz inequality for $\mathbb{E}_{\mathbb{Q}} [|\zeta_s(y)| \cdot \psi(y, x)^2] \leq \mathbb{E}_{\mathbb{Q}} [|\zeta_s(y)|^2]^{1/2} \cdot$
 1605 $\mathbb{E}_{\mathbb{Q}} [\psi(y, x)^4]^{1/2} \leq \mathbb{E}_{\mathbb{Q}} [\psi(y, x)^4]^{1/2} = O(1)$. As for the case $s^* = 1$, we already have a constant
 1606 outside, and noting that the second moment is at most $O(1)$ due to the quadruple-integrable as-
 1607 sumption, it suffices to consider in the following $s^* \geq 2$. For the second part of Eq. (H.8), we
 1608 can split the expectation according to Eq. (H.7). For the first term in Eq. (H.7) which corresponds
 1609 to $\sqrt{(s+2)(s+1)} \cdot \mathbf{h}_{s+2}(z) [(\theta^*)^{\otimes s} \otimes v^{\otimes 2}]$, we take test tensor $T_s = v^{\otimes 2} \otimes (\theta^*)^{\otimes (s-2)}$ with
 1610 $c_0 = 2, s_0 = s^* + 2$ and have by Proposition I.4 that

$$\begin{aligned} & \left| \sum_{s=s^*}^{\infty} \sqrt{(s+2)(s+1)} \cdot \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \mathbf{h}_{s+2}(z) [(\theta^*)^{\otimes s} \otimes v^{\otimes 2}]] \right| \\ &\lesssim \mathbb{1}(s_0 \leq c_0) \cdot \left(\epsilon_1^{s^*-1-\lfloor s_0/2 \rfloor} \cdot \epsilon_0^{s_0} + \epsilon_1^{s^*-1-\lfloor c_0/2 \rfloor} \cdot \epsilon_0^{c_0} \right) + \epsilon_0^{c_0} \cdot \epsilon^{(2s^*) \vee s_0 - c_0} \\ &\quad + \mathbb{1}(s_0 \leq 2(s^* - 1)) \cdot \left(\epsilon_1^{s^*-1-\lfloor (c_0+1)/2 \rfloor} \cdot \epsilon_0^{c_0} \cdot \epsilon + \epsilon_0^{c_0} \cdot \epsilon^{2(s^*-1)+1-c_0} \right) \\ &\lesssim \epsilon_0^2 \cdot \epsilon^{2s^*-2} + \mathbb{1}(s^* \geq 4) \cdot \left(\epsilon_1^{s^*-2} \cdot \epsilon_0^2 \cdot \epsilon + \epsilon_0^2 \cdot \epsilon^{2s^*-3} \right). \end{aligned}$$

1611 For the second term $\mathbf{h}_s(z) [(\theta^*)^{\otimes s}]$, we take test tensor $T_s = (\theta^*)^{\otimes s}$ with $c_0 = 0, s_0 = s^*$ and have
 1612 by Proposition I.4 that

$$\begin{aligned} & \left| \sum_{s=s^*}^{\infty} \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \mathbf{h}_s(z) [(\theta^*)^{\otimes s}]] \right| \\ &\lesssim \epsilon_0^{c_0} \cdot \epsilon^{(2s^*) \vee s_0 - c_0} + \left(\epsilon_1^{s^*-1-\lfloor (c_0+1)/2 \rfloor} \cdot \epsilon_0^{c_0} \cdot \epsilon + \epsilon_0^{c_0} \cdot \epsilon^{2(s^*-1)+1-c_0} \right) \\ &\lesssim \epsilon^{2s^*} + \epsilon_1^{s^*-1} \cdot \epsilon + \epsilon^{2s^*-1} \lesssim \epsilon_1^{s^*-1} \cdot \epsilon + \epsilon^{2s^*-1}. \end{aligned}$$

1613 For the third term $2s \cdot \mathbf{h}_s(z) [(\theta^*)^{\otimes s-1} \otimes v]^{\top} \cdot \langle \theta^*, v \rangle$, we take test tensor $T_s = v \otimes (\theta^*)^{\otimes s-1}$ with
 1614 $c_0 = 1, s_0 = s^*$ and have by Proposition I.4 that

$$\begin{aligned} & \left| \sum_{s=s^*}^{\infty} 2s \cdot \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \mathbf{h}_s(z) [(\theta^*)^{\otimes s-1} \otimes v]^{\top}] \right| \\ &\lesssim \epsilon_0^{c_0} \cdot \epsilon^{(2s^*) \vee s_0 - c_0} + \left(\epsilon_1^{s^*-1-\lfloor (c_0+1)/2 \rfloor} \cdot \epsilon_0^{c_0} \cdot \epsilon + \epsilon_0^{c_0} \cdot \epsilon^{2(s^*-1)+1-c_0} \right) \\ &\lesssim \epsilon_0 \cdot \epsilon^{2s^*-1} + \epsilon_1^{s^*-2} \cdot \epsilon_0 \cdot \epsilon + \epsilon_0 \cdot \epsilon^{2s^*-2} \lesssim \epsilon_1^{s^*-2} \cdot \epsilon_0 \cdot \epsilon + \epsilon_0 \cdot \epsilon^{2s^*-2}. \end{aligned}$$

1615 For the last term $\sqrt{s(s-1)} \cdot \mathbf{h}_{s-2}(z) [(\theta^*)^{\otimes s-2}] \cdot \langle \theta^*, v \rangle^2$, we take test tensor $T_s = (\theta^*)^{\otimes s}$ with
 1616 $c_0 = 0, s_0 = s^* - 2$ and have by Proposition I.4 that

$$\begin{aligned} & \left| \sum_{s=s^*}^{\infty} \sqrt{s(s-1)} \cdot \mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \mathbf{h}_{s-2}(z) [(\theta^*)^{\otimes s-2}] \cdot \langle \theta^*, v \rangle^2] \right| \\ &\lesssim \mathbb{1}(s^* = 2) \cdot \left(\epsilon_1^{s^*-1-\lfloor s_0/2 \rfloor} \cdot \epsilon_0^{s_0} + \epsilon_1^{s^*-1-\lfloor c_0/2 \rfloor} \cdot \epsilon_0^{c_0} \right) + \epsilon_0^{c_0} \cdot \epsilon^{(2s^*) \vee s_0 - c_0} \\ &\quad + \left(\epsilon_1^{s^*-1-\lfloor (c_0+1)/2 \rfloor} \cdot \epsilon_0^{c_0} \cdot \epsilon + \epsilon_0^{c_0} \cdot \epsilon^{2(s^*-1)+1-c_0} \right) \\ &\lesssim \mathbb{1}(s^* = 2) \epsilon_1 + \epsilon^{2s^*} + \epsilon_1^{s^*-1} \cdot \epsilon + \epsilon^{2s^*-1} \lesssim \mathbb{1}(s^* = 2) \epsilon_1 + \epsilon_1^{s^*-1} \cdot \epsilon + \epsilon^{2s^*-1}. \end{aligned}$$

1617 Summing up the above terms, we have that for $s^* \geq 2$,

$$\begin{aligned} F &\lesssim \left(\epsilon_1^{s^*-2} \epsilon_0^2 + \epsilon_1^{s^*-1} \right) + \left(\epsilon_0^2 \cdot \epsilon^{2s^*-2} + \mathbb{1}(s^* \geq 4) \cdot \left(\epsilon_1^{s^*-2} \cdot \epsilon_0^2 \cdot \epsilon + \epsilon_0^2 \cdot \epsilon^{2s^*-3} \right) \right) \\ &\quad + \left(\epsilon_1^{s^*-1} \cdot \epsilon + \epsilon^{2s^*-1} \right) + \left(\epsilon_1^{s^*-2} \cdot \epsilon_0 \cdot \epsilon + \epsilon_0 \cdot \epsilon^{2s^*-2} + \mathbb{1}(s^* = 2) \epsilon_1 + \epsilon_1^{s^*-1} \cdot \epsilon + \epsilon^{2s^*-1} \right) \mathbb{1}(v = \theta^*) \\ &\lesssim \epsilon_1^{s^*-2} \epsilon_0^2 + \epsilon_1^{s^*-1} + \epsilon_0^2 \cdot \epsilon^{2s^*-2} + \mathbb{1}(s^* \geq 4) \cdot \epsilon_0^2 \cdot \epsilon^{2s^*-3} + \epsilon^{2s^*-1} \\ &\quad + \left(\epsilon_1^{s^*-2} \cdot \epsilon_0 \cdot \epsilon + \epsilon_0 \cdot \epsilon^{2s^*-2} \right) \mathbb{1}(v = \theta^*). \end{aligned}$$

1618 If $v = \theta^*$, then we additionally have $\epsilon_0 = \epsilon$, which simplifies the above bound to

$$F|_{s^* \geq 2, v = \theta^*} \lesssim \epsilon_1^{s^*-2} \epsilon^2 + \epsilon_1^{s^*-1} + \epsilon^{2s^*-1} = \epsilon_1^{s^*-1} \cdot \left(1 + \frac{\epsilon^2}{\epsilon_1} + \left(\frac{\epsilon^2}{\epsilon_1} \right)^{s^*-1} \cdot \epsilon \right).$$

1619 For $v \perp \theta^*$, we have that

$$\begin{aligned} F|_{s^* \geq 2, v \perp \theta^*} &\lesssim \epsilon_1^{s^*-2} \epsilon_0^2 + \epsilon_1^{s^*-1} + \epsilon_0^2 \cdot \epsilon^{2s^*-2} + \mathbb{1}(s^* \geq 4) \cdot \epsilon_0^2 \cdot \epsilon^{2s^*-3} + \epsilon^{2s^*-1} \\ &\lesssim \epsilon_1^{s^*-1} \cdot \left(1 + \frac{\epsilon^2}{\epsilon_1} + \left(\frac{\epsilon^2}{\epsilon_1} \right)^{s^*-1} \cdot \epsilon + \left(\frac{\epsilon^2}{\epsilon_1} \right)^{s^*-2} \cdot \frac{\epsilon_0^2}{\epsilon_1} \cdot (\epsilon^2 + \epsilon \cdot \mathbb{1}(s^* \geq 4)) \right). \end{aligned}$$

1620 Where for $s^* = 1$, we have $F|_{s^*=1} \lesssim 1$. Hence, we complete the proof. \square

1621 **Lemma H.4.** For polarization level $\gamma = o(1)$, take the polarized random vector

$$w = \frac{\gamma e_1 + \xi}{\|\gamma e_1 + \xi\|_2}, \quad \text{where } \xi \sim \text{Unif}(\mathbb{S}^{d-1}),$$

1622 where $e_1 = (1, 0, \dots, 0)^\top$ is the first standard basis vector in \mathbb{R}^d . Let $\theta^* = (\rho, \sqrt{1-\rho^2}, 0, \dots, 0) \in$
1623 \mathbb{S}^{d-1} be a fixed direction. Then, we have that

$$\mathbb{E}[\langle \theta^*, w \rangle^s] \simeq \begin{cases} \left(|\rho|(\gamma + d^{-1/2}) + \sqrt{1-\rho^2} d^{-1/2} \right)^s & \text{if } s \text{ is even} \\ \rho \gamma \left(|\rho|(\gamma + d^{-1/2}) + \sqrt{1-\rho^2} d^{-1/2} \right)^{s-1} & \text{if } s \text{ is odd.} \end{cases}$$

1624 *Proof of Lemma H.4.* For w , we have by Lemma I.9 that the first moment is given by

$$\mathbb{E}[w] = C(e_1, \gamma) \cdot \gamma \cdot e_1 \simeq \gamma e_1,$$

1625 and the second moment is controlled by

$$\mathbb{E}[w w^\top] = \begin{bmatrix} C(2e_1, \gamma) \cdot (\gamma + d^{-1/2})^2 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & C(2e_2, \gamma) d^{-1} \cdot I_{d-1} & & \\ 0 & & & \end{bmatrix} \lesssim \begin{bmatrix} (\gamma + d^{-1/2})^2 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & & \\ 0 & & & d^{-1} \cdot I_{d-1} \end{bmatrix}.$$

1626 For $\langle \theta^*, w \rangle^s \cdot w$, we look at coordinate τ of w , and the first moment is given by

$$\mathbb{E}[\langle \theta^*, w \rangle^s \cdot w_\tau] = \mathbb{E} \left[(\rho w_1 + \sqrt{1-\rho^2} w_2)^s w_\tau \right].$$

1627 Note that if $\tau \neq 1, 2$, then the expectation is zero. For more generality, let us take $r_1, r_2 \in \mathbb{N}$. We
1628 study the following expectation:

$$\begin{aligned} &\mathbb{E}[\langle \theta^*, w \rangle^s \cdot w_1^{r_1} w_2^{r_2}] \\ &= \sum_{j=0}^s \binom{s}{j} \rho^j \sqrt{1-\rho^2}^{s-j} \mathbb{E} \left[w_1^{j+r_1} w_2^{s-j+r_2} \right] \cdot \mathbb{1}(s-j+r_2 \text{ even}) \\ &= \begin{cases} \sum_{j=0}^{\lfloor s/2 \rfloor} \binom{s}{2j} \rho^{2j} \sqrt{1-\rho^2}^{s-2j} \mathbb{E} \left[w_1^{2j+r_1} w_2^{s-2j+r_2} \right] & \text{if } s+r_2 \text{ is even,} \\ \sum_{j=0}^{\lfloor (s-1)/2 \rfloor} \binom{s}{2j+1} \rho^{2j+1} \sqrt{1-\rho^2}^{s-2j-1} \mathbb{E} \left[w_1^{2j+1+r_1} w_2^{s-2j-1+r_2} \right] & \text{if } s+r_2 \text{ is odd.} \end{cases} \end{aligned}$$

1629 Here, the first equality holds by noting that each term in the sum is zero if the degree on w_2 is odd
 1630 due to the symmetry in the distribution of w_2 . Next, we invoke Lemma I.9 for the moment as

$$\begin{aligned} & \mathbb{E}[\langle \theta^*, w \rangle^s \cdot w_1^{r_1} w_2^{r_2}] \\ & \simeq \begin{cases} \text{if } s + r_2 \text{ is even:} \\ \sum_{j=0}^{\lfloor s/2 \rfloor} \binom{s}{2j} \rho^{2j} \sqrt{1 - \rho^2}^{s-2j} \gamma^{\mathbb{1}(r_1 \text{ odd})} (\gamma + d^{-1/2})^{2j+2\lfloor r_1/2 \rfloor} (d^{-1/2})^{s-2j+r_2} \\ \text{if } s + r_2 \text{ is odd:} \\ \sum_{j=0}^{\lfloor (s-1)/2 \rfloor} \binom{s}{2j+1} \rho^{2j+1} \sqrt{1 - \rho^2}^{s-2j-1} \gamma^{\mathbb{1}(r_1 \text{ even})} (\gamma + d^{-1/2})^{2j+2\lfloor \frac{r_1+1}{2} \rfloor} (d^{-1/2})^{s-2j-1+r_2} \end{cases} \\ & \simeq \begin{cases} \text{if } s + r_2 \text{ is even:} \\ \sqrt{1 - \rho^2}^{\mathbb{1}(s \text{ odd})} \gamma^{\mathbb{1}(r_1 \text{ odd})} \left(\gamma + \frac{1}{\sqrt{d}} \right)^{2\lfloor \frac{r_1}{2} \rfloor} \left(\frac{1}{\sqrt{d}} \right)^{r_2 + \mathbb{1}(s \text{ odd})} \left(|\rho|(\gamma + d^{-1/2}) + \sqrt{1 - \rho^2} d^{-1/2} \right)^{2\lfloor \frac{s}{2} \rfloor} \\ \text{if } s + r_2 \text{ is odd:} \\ \rho \sqrt{1 - \rho^2}^{\mathbb{1}(s \text{ even})} \gamma^{\mathbb{1}(r_1 \text{ even})} \left(\gamma + \frac{1}{\sqrt{d}} \right)^{2\lfloor \frac{r_1+1}{2} \rfloor} \left(\frac{1}{\sqrt{d}} \right)^{r_2 + \mathbb{1}(s \text{ even})} \left(|\rho|(\gamma + d^{-1/2}) + \sqrt{1 - \rho^2} d^{-1/2} \right)^{2\lfloor \frac{s-1}{2} \rfloor} \end{cases} \end{aligned}$$

1631 Here, the symbol \simeq conceals some constant factors that are governed by upper and lower bounds
 1632 dependent on s only. Using the above calculation, we have We can specialize the above results to the
 1633 case $r_1 = r_2 = 0$ and obtain

$$\mathbb{E}[\langle \theta^*, w \rangle^s] \simeq \begin{cases} \left(|\rho|(\gamma + d^{-1/2}) + \sqrt{1 - \rho^2} d^{-1/2} \right)^s & \text{if } s \text{ is even} \\ \rho \gamma \left(|\rho|(\gamma + d^{-1/2}) + \sqrt{1 - \rho^2} d^{-1/2} \right)^{s-1} & \text{if } s \text{ is odd,} \end{cases}$$

1634 which completes the proof. \square

1635 I Technical Results

1636 I.1 Technical Results for Hermite Tensor

1637 **Proposition I.1.** *Let $s \in \mathbb{N}_0$. For any $z \in \mathbb{R}^d$, we have*

$$\begin{aligned} z_i z_j \mathbf{h}_s(z) &= \text{Sym}(\sqrt{(s+2)(s+1)} \cdot \mathbf{h}_{s+2}(z)[e_i \otimes e_j] + \delta_{ij} \mathbf{h}_s(z) + s \cdot \mathbf{h}_s(z)[e_j] \otimes e_i \\ & \quad + s \cdot \mathbf{h}_s(z)[e_i] \otimes e_j + \sqrt{s(s-1)} \cdot \mathbf{h}_{s-2}(z) \otimes e_i \otimes e_j), \end{aligned}$$

1638 where we define $\mathbf{h}_{-1}(z)$ and $\mathbf{h}_{-2}(z)$ to be all zero tensors of any conformable shape.

1639 *Proof of Proposition I.1.* Note that each element of $\mathbf{h}_s(\theta^{*\top} z) z z^\top$ must lie in the polynomial space
 1640 with degree at most $s+2$, i.e., $\mathbb{R}_{s+2}[z]$. We take a test function $F: \mathbb{R}^d \rightarrow \mathbb{R}$ such that $R \in \mathbb{R}_{s+2}[z]$.
 1641 Thus, we can write down the inner product of F and $\mathbf{h}_s(\theta^{*\top} z) z_i z_j$ for $i, j \in [d]$ as

$$\mathbb{E}_{z \sim \mathcal{N}_d}[F(z) \mathbf{h}_s(\theta^{*\top} z) z_i z_j] = \mathbb{E}_{z \sim \mathcal{N}_d}[F(z) z_i z_j \cdot \mathbf{h}_s(z) [(\theta^*)^{\otimes s}]],$$

1642 where in the equation, we use (B.4) to rewrite the Hermite polynomial in terms of the Hermite tensor.
 1643 It suffices to understand the tensor $F(z) z_i z_j \mathbf{h}_s(z)$. Note that F is differentiable to any order, and the
 1644 tensor obtained by differentiating F to any order is square-integrable with respect to the standard
 1645 normal distribution. By the Stein's lemma for Hermite tensor (B.5), we have

$$\begin{aligned} \sqrt{s!} \cdot \mathbb{E}_{z \sim \mathcal{N}_d}[F(z) z_i z_j \mathbf{h}_s(z)] &= \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s(F(z) z_i z_j)] \\ &= \text{Sym} \left(\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) z_i z_j] + s \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s-1} F(z) \nabla(z_i z_j)] + \frac{s(s-1)}{2} \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s-2} F(z) \nabla^2(z_i z_j)] \right) \\ &= \text{Sym} \left(\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) \delta_{ij}] + \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s+2} F(z) [e_i \otimes e_j]] + s \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) [e_j] \otimes e_i] \right. \\ & \quad \left. + s \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) [e_i] \otimes e_j] + s(s-1) \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s-2} F(z) \otimes e_i \otimes e_j] \right). \end{aligned} \tag{I.1}$$

1646 Here, the ‘‘Sym’’ operation symmetrizes the tensor in the parentheses. The last equality holds by the
 1647 following calculations. For $\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) z_i z_j]$, we have

$$\begin{aligned} \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) z_i z_j] &= \text{Sym} \left(\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) \delta_{ij}] + \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s+1} F(z)[e_j] z_i] \right) \\ &= \text{Sym} \left(\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z) \delta_{ij}] + \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s+2} F(z)[e_i \otimes e_j]] \right), \end{aligned}$$

1648 where we use the Stein’s lemma for both equalities. For $\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s-2} F(z) \nabla^2(z_i z_j)]$, we have

$$\begin{aligned} \text{Sym} \left(\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^{s-1} F(z) \nabla(z_i z_j)] \right) &= \text{Sym} \left(\mathbb{E}_{z \sim \mathcal{N}_d}[z_j \nabla^{s-1} F(z) \otimes e_i + z_i \nabla^{s-1} F(z) \otimes e_j] \right) \\ &= \text{Sym} \left(\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z)[e_j] \otimes e_i] + \mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z)[e_i] \otimes e_j] \right). \end{aligned}$$

1649 Now, for each derivative of F in (I.1), we have by the Stein’s lemma stated in (B.5) that
 1650 $\mathbb{E}_{z \sim \mathcal{N}_d}[\nabla^s F(z)] = \sqrt{s!} \cdot \mathbb{E}_{z \sim \mathcal{N}_d}[F(z) \mathbf{h}_s(z)]$, which gives us

$$\begin{aligned} \mathbb{E}_{z \sim \mathcal{N}_d}[F(z) z_i z_j \mathbf{h}_s(z)] &= \mathbb{E}_{z \sim \mathcal{N}_d} \left[F(z) \cdot \text{Sym}(\delta_{ij} \mathbf{h}_s(z) + \sqrt{(s+2)(s+1)} \cdot \mathbf{h}_{s+2}(z)[e_i \otimes e_j] \right. \\ &\quad \left. + s \cdot \mathbf{h}_s(z)[e_j] \otimes e_i + s \cdot \mathbf{h}_s(z)[e_i] \otimes e_j + \sqrt{s(s-1)} \cdot \mathbf{h}_{s-2}(z) \otimes e_i \otimes e_j \right]. \end{aligned}$$

1651 Since $F \in \mathbb{R}_{s+2}[z]$ is arbitrary, we conclude that

$$\begin{aligned} z_i z_j \mathbf{h}_s(z) &= \text{Sym} \left(\sqrt{(s+2)(s+1)} \cdot \mathbf{h}_{s+2}(z)[e_i \otimes e_j] + \delta_{ij} \mathbf{h}_s(z) + s \cdot \mathbf{h}_s(z)[e_j] \otimes e_i \right. \\ &\quad \left. + s \cdot \mathbf{h}_s(z)[e_i] \otimes e_j + \sqrt{s(s-1)} \cdot \mathbf{h}_{s-2}(z) \otimes e_i \otimes e_j \right). \end{aligned}$$

1652 The proof is completed by further taking the tensor inner product operation with respect to $(\theta^*)^{\otimes s}$ on
 1653 both side. \square

1654 **Proposition I.2.** Let $w, w' \in \mathbb{S}^{d-1}$ and $s \in \mathbb{N}_0$. We have

$$\begin{aligned} &\mathbb{E} [h_i(\langle w, z \rangle) h_j(\langle w', z \rangle) \cdot \mathbf{h}_s(z)] \\ &= \sum_{\tau=0}^s \mathbb{1}(j = i + s - 2\tau, i \geq \tau) \cdot \binom{s}{\tau} \sqrt{\frac{i!j!}{s!((i-\tau)!)^2}} \cdot \langle w, w' \rangle^{i-\tau} \cdot \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes s-\tau}). \end{aligned}$$

1655 The above term is equal to

$$\binom{s}{(i-j+s)/2} \cdot \sqrt{\frac{i!j!}{s!}} \cdot \frac{\langle w, w' \rangle^{(i+j-s)/2}}{((i+j-s)/2)!} \cdot \text{Sym}(w^{\otimes (i-j+s)/2} \otimes w'^{\otimes (j-i+s)/2}).$$

1656 if $|i-j| \leq s \leq i+j$ and $s \equiv i-j \pmod{2}$. Otherwise the expectation gives the zero tensor.

1657 *Proof of Proposition I.2.* By the Stein’s lemma for the Hermite tensor (B.5), we have

$$\begin{aligned} \mathbb{E} [h_i(\langle w, z \rangle) h_j(\langle w', z \rangle) \cdot \mathbf{h}_s(z)] &= \frac{1}{\sqrt{s!}} \cdot \mathbb{E} [\nabla^s (h_i(\langle w, z \rangle) h_j(\langle w', z \rangle))] \\ &= \frac{1}{\sqrt{s!}} \cdot \sum_{\tau=0}^s \binom{s}{\tau} \mathbb{E} [\text{Sym}(\nabla^\tau h_i(\langle w, z \rangle) \otimes \nabla^{s-\tau} h_j(\langle w', z \rangle))] \\ &= \frac{1}{\sqrt{s!}} \cdot \sum_{\tau=0}^s \binom{s}{\tau} \sqrt{\frac{i!j!}{(i-\tau)!(j-s+\tau)!}} \mathbb{E} [h_{i-\tau}(\langle w, z \rangle) \otimes h_{j-s+\tau}(\langle w', z \rangle) \cdot \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes s-\tau})] \\ &= \sum_{\tau=0}^s \mathbb{1}(j = i + s - 2\tau, i \geq \tau) \cdot \frac{\binom{s}{\tau}}{(i-\tau)!} \sqrt{\frac{i!j!}{s!}} \langle w, w' \rangle^{i-\tau} \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes s-\tau}). \end{aligned}$$

1658 the condition can be translated into $|i-j| \leq s \leq i+j$ and $s \equiv i-j \pmod{2}$. Then, we can take
 1659 $\tau = (i-j+s)/2, s-\tau = (j-i+s)/2, i-\tau = (i+j-s)/2$ to obtain the desired result. \square

1660 **Lemma I.3.** Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ such that $\psi^2 \in L^2(\mathcal{N})$. Suppose that ψ is high-pass in the sense
 1661 that $\hat{\psi}_i = 0$ for any $i < s^* - 1$ for some $s^* \in \mathbb{N}_0$. For $w, w' \in \mathbb{S}^{d-1}$, take a series of test tensor
 1662 $\{T_s = v_1 \otimes v_2 \otimes \dots \otimes v_s \in (\mathbb{R}^d)^{\otimes s}\}_{s=0}^\infty$ such that $\sup_{i > c_0} \{|\langle w, v_i \rangle| \vee |\langle w', v_i \rangle|\} \leq \epsilon$ for some
 1663 $\epsilon \in (0, 1/2)$ and integer $c_0 \in \mathbb{N}_0$. Let

$$c_0 := \max_{1 \leq i \leq c_0} \{|\langle w, v_i \rangle| \vee |\langle w', v_i \rangle|\}.$$

1664 Suppose that $|\langle w, w' \rangle| \leq \epsilon_1$. Then we have for any $s \in \mathbb{N}_0$ that

$$\left| \mathbb{E}_{z \sim \mathcal{N}_d} [\psi(\langle w, z \rangle) \psi(\langle w', z \rangle) \cdot \mathbf{h}_s(z) [T_s]] \right| \leq 4 \|\psi\|_2^2 (4e s^*)^{s/2} \sqrt{s + s^*} \cdot \epsilon_1^{(s^* - 1 - \lfloor s/2 \rfloor) \vee 0} \cdot \epsilon_0^{s \wedge c_0} \cdot \epsilon^{(s - c_0) \vee 0},$$

1665 where $\|\psi\|_2^2 = \mathbb{E}_{x \sim \mathcal{N}} [\psi^2(x)]$.

1666 *Proof of Lemma I.3.* As $\psi(x)\psi(x')$ is also square-integrable, we are able to extract the s -th tensor
1667 coefficient of the Hermite expansion of $\psi(x)\psi(x')$ as

$$\begin{aligned} & \mathbb{E}_{z \sim \mathcal{N}_d} [\psi(\langle w, z \rangle) \psi(\langle w', z \rangle) \cdot \mathbf{h}_s(z)] \\ &= \mathbb{E}_{z \sim \mathcal{N}_d} \left[\sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \hat{\psi}_i \hat{\psi}_j h_i(\langle w, z \rangle) h_j(\langle w', z \rangle) \cdot \mathbf{h}_s(z) \right] \\ &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \sum_{\tau=0}^s \mathbb{1}(j = i + s - 2\tau, i \geq \tau) \cdot \binom{s}{\tau} \sqrt{\frac{i!j!}{s!((i-\tau)!)^2}} \hat{\psi}_i \hat{\psi}_j \cdot \langle w, w' \rangle^{i-\tau} \cdot \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes s-\tau}) \\ &= \sum_{i=0}^{\infty} \sum_{\tau=0}^s \mathbb{1}(i \geq \tau) \binom{s}{\tau} \sqrt{\frac{1}{s!}} \sqrt{\frac{i!(i+s-2\tau)!}{((i-\tau)!)^2}} \hat{\psi}_i \hat{\psi}_{i+s-2\tau} \cdot \langle w, w' \rangle^{i-\tau} \cdot \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes (s-\tau)}), \end{aligned}$$

1668 where the last second identity follows from Proposition I.2, and in the last line we restrict the condition
1669 $j = i + s - 2\tau$. Note that the double sums are interchangeable only if the series converges for each τ .
1670 However, by our condition that $\langle w, w' \rangle \leq 1 - \epsilon$ for some $\epsilon > 0$, then for each τ , we have for any test
1671 tensor $T = v_1 \otimes v_2 \otimes \dots \otimes v_s$ with $\|v_i\|_2 = 1$ that

$$\begin{aligned} & \left| \sum_{i=\tau}^{\infty} \sqrt{\frac{i!(i+s-2\tau)!}{((i-\tau)!)^2}} \hat{\psi}_i \hat{\psi}_{i+s-2\tau} \cdot \langle w, w' \rangle^{i-\tau} \cdot \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes (s-\tau)}) [T] \right| \\ & \leq \sum_{i=\tau}^{\infty} (i+s)^{s/2} \cdot \left| \hat{\psi}_i \hat{\psi}_{i+s-2\tau} \right| \cdot (1-\epsilon)^{i-\tau} < \infty, \end{aligned}$$

1672 where we note that $1 - \epsilon$ will dominate the polynomial growth of $(i+s)^{s/2}$, and also using the fact
1673 that $\hat{\psi}_i \hat{\psi}_{i+s-2\tau}$ is uniformly bounded by $\sum_{i=0}^{\infty} \hat{\psi}_i^2 < \infty$. Now, we interchange the double sum and
1674 apply the high-pass assumption and have that

$$\begin{aligned} & \mathbb{E}_{z \sim \mathcal{N}_d} [\psi(\langle w, z \rangle) \psi(\langle w', z \rangle) \cdot \mathbf{h}_s(z) [T]] \\ &= \sum_{\tau=0}^s \binom{s}{\tau} \sqrt{\frac{1}{s!}} \underbrace{\sum_{i=i_0}^{\infty} \sqrt{\frac{i!(i+s-2\tau)!}{((i-\tau)!)^2}} \hat{\psi}_i \hat{\psi}_{i+s-2\tau} \cdot \langle w, w' \rangle^{i-\tau} \cdot \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes (s-\tau)}) [T]}_{\text{(I)}}. \end{aligned}$$

1675 where we define

$$i_0 = \max \{ (s^* - 1), (s^* - 1 + 2\tau - s), \tau \}$$

1676 Note that the tensor product part is independent of i . Hence, we pull out term (I) and have

$$\begin{aligned} \text{(I)} &= \sqrt{\frac{i_0!(i_0+s-2\tau)!}{((i_0-\tau)!)^2}} \hat{\psi}_{i_0} \hat{\psi}_{i_0+s-2\tau} \langle w, w' \rangle^{i_0-\tau} \\ & \pm \|\psi\|_2^2 \cdot \langle w, w' \rangle^{i_0-\tau+1} \left(\sum_{j=0}^{\infty} (j+s+s^*)(j+s+s^*-1) \cdots (j+s^* + \lceil (s+1)/2 \rceil) |\langle w, w' \rangle|^j \right). \end{aligned}$$

1677 Here, the first term is given by splitting out the term with $i = i_0$ from the summation, and the second
1678 term is for $i \geq i_0 + 1$. For the second term, we have the following argument:

$$\begin{aligned} \max\{i_0, i_0 + s - 2\tau\} &= (s^* - 1) \vee (s^* - 1 + 2\tau - s) \vee \tau \vee (s^* - 1 + s - 2\tau) \vee (s^* - 1) \vee (s - \tau) \\ &= (s^* - 1 + 2\tau - s) \vee (s^* - 1 + s - 2\tau) \vee \tau \vee (s - \tau) \\ &= (s^* - 1 + |2\tau - s|) \vee (|\tau - s/2| + s/2) \leq s^* + s - 1. \end{aligned}$$

1679 Hence, we then have that for any $i \geq i_0 + 1$ with $j = i - (i_0 + 1)$ that

$$\max\{i, i + s - 2\tau\} \leq s^* + s + j.$$

1680 Therefore, for any $i \geq i_0 + 1$, we have

$$\begin{aligned} \sqrt{\frac{i!(i+s-2\tau)!}{((i-\tau)!)^2}} &= \sqrt{i(i-1)\cdots(i-\tau+1)} \cdot \sqrt{(i+s-2\tau)(i+s-2\tau-1)\cdots(i-\tau+1)} \\ &\leq \sqrt{(j+s^*+s)(j+s^*+s-1)\cdots(j+s^*+1)} \\ &\leq (j+s^*+s)(j+s^*+s-1)\cdots(j+s^*+\lceil(s+1)/2\rceil). \end{aligned}$$

1681 To characterize the second term, we use Proposition J.5 where we have conditions $|\langle w, w' \rangle| < 1/2$
1682 and $s^* + s \geq 2\lceil(s+1)/2\rceil - 1$ satisfied, which gives us

$$\begin{aligned} &\sum_{j=0}^{\infty} (j+s+s^*)(j+s+s^*-1)\cdots(j+s^*+\lceil(s+1)/2\rceil) |\langle w, w' \rangle|^j \\ &\leq 2(s+s^*)\cdots(s^*+\lceil(s+1)/2\rceil) \cdot \frac{1}{1-|\langle w, w' \rangle|} \\ &\leq 4(s+s^*)\cdots(s^*+\lceil(s+1)/2\rceil). \end{aligned}$$

1683 Combining these results, we have for (I) that

$$\begin{aligned} \text{(I)} &= \sqrt{\frac{i_0!(i_0+s-2\tau)!}{((i_0-\tau)!)^2}} \widehat{\psi}_{i_0} \widehat{\psi}_{i_0+s-2\tau} \langle w, w' \rangle^{i_0-\tau} \pm \|\psi\|_2^2 \cdot 4(s+s^*)\cdots(s^*+\lceil(s+1)/2\rceil) \cdot \langle w, w' \rangle^{i_0-\tau+1} \\ &= C(s^*, s, \tau, \langle w, w' \rangle) \cdot \langle w, w' \rangle^{i_0-\tau}, \end{aligned}$$

1684 where we define $C(s^*, s, \tau, \langle w, w' \rangle) = \text{(I)}/\langle w, w' \rangle^{i_0-\tau}$ as the coefficient, which is given by

$$C(s^*, s, \tau, \langle w, w' \rangle) = \sqrt{\frac{i_0!(i_0+s-2\tau)!}{((i_0-\tau)!)^2}} \widehat{\psi}_{i_0} \widehat{\psi}_{i_0+s-2\tau} \pm \|\psi\|_2^2 \cdot 4(s+s^*)\cdots(s^*+\lceil(s+1)/2\rceil) \cdot |\langle w, w' \rangle|,$$

1685 and also enjoys the following upper bound

$$|C(s^*, s, \tau, \langle w, w' \rangle)| \leq \|\psi\|_2^2 \cdot 4(s+s^*-1)\cdots(s^*-1+\lceil(s+1)/2\rceil). \quad (\text{I.2})$$

1686 Here, the upper bound can be obtained by noting that the previous upper bound for terms $i \geq i_0 + 1$
1687 can be also applied to $i \geq i_0$.

1688 **Case $s \geq 1$.** Let us plug in test tensor T_s into the expression, which gives us

$$\begin{aligned} &|\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(\langle w, z \rangle) \psi(\langle w', z \rangle) \cdot \mathbf{h}_s(z) [T_s]]| \\ &\leq \sum_{\tau=0}^s \binom{s}{\tau} \sqrt{\frac{1}{s!}} |C(s^*, s, \tau, \langle w, w' \rangle)| \cdot |\langle w, w' \rangle|^{i_0-\tau} \cdot \left| \text{Sym}(w^{\otimes \tau} \otimes w'^{\otimes (s-\tau)}) [T_s] \right| \\ &\leq \sum_{\tau=0}^s \binom{s}{\tau} \sqrt{\frac{1}{s!}} |C(s^*, s, \tau, \langle w, w' \rangle)| \cdot |\langle w, w' \rangle|^{i_0-\tau} \cdot \left(\frac{1}{s!} \sum_{\pi \in \Pi_s} \prod_{i=1}^{\tau} |\langle w, v_{\pi(i)} \rangle| \prod_{j=\tau+1}^s |\langle w', v_{\pi(j)} \rangle| \right), \end{aligned}$$

1689 where Π_s denotes the set of all permutations of s elements. We further have this term bounded by

$$\begin{aligned} &|\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(\langle w, z \rangle) \psi(\langle w', z \rangle) \cdot \mathbf{h}_s(z) [T_s]]| \quad (\text{I.3}) \\ &\leq \sum_{\tau=0}^s \binom{s}{\tau} \sqrt{\frac{1}{s!}} \max_{0 \leq \tau \leq s} |C(s^*, s, \tau, \langle w, w' \rangle)| \cdot |\langle w, w' \rangle|^{i_0-\tau} \cdot \epsilon_0^{s \wedge c_0} \cdot \epsilon^{(s-c_0) \vee 0} \\ &\leq \frac{2^s}{\sqrt{s!}} \max_{0 \leq \tau \leq s} |C(s^*, s, \tau, \langle w, w' \rangle)| \cdot |\langle w, w' \rangle|^{(s^*-1-\lfloor s/2 \rfloor) \vee 0} \cdot \epsilon_0^{s \wedge c_0} \cdot \epsilon^{(s-c_0) \vee 0}. \end{aligned}$$

1690 where the last inequality follows from the following fact

$$\begin{aligned} i_0 - \tau &= (s^* - 1 - \tau) \vee (s^* - 1 + \tau - s) \vee 0 = (s^* - 1 - s/2 + |\tau - s/2|) \vee 0 \\ &\geq (s^* - 1 - s/2 + \mathbf{1}(s \text{ odd})/2) \vee 0 = (s^* - 1 - \lfloor s/2 \rfloor) \vee 0. \end{aligned}$$

1691 **Plugging (I.2) into (I.3), and by noting that $|\langle w, v_i \rangle| \vee |\langle w', v_i \rangle| \leq \epsilon$ for any $v_i \in \{\theta^*, v\}$ and**
 1692 **$|\langle w, w' \rangle| \leq \epsilon$, we have that**

$$\begin{aligned} & |\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(\langle w, z \rangle) \psi(\langle w', z \rangle) \cdot \mathbf{h}_s(z) [T_s]]| \\ & \leq \frac{2^s}{\sqrt{s!}} \|\psi\|_2^2 \cdot 4(s + s^* - 1) \cdots (s^* - 1 + \lceil (s+1)/2 \rceil) \cdot \epsilon_1^{(s^* - 1 - \lfloor s/2 \rfloor) \vee 0} \cdot \epsilon_0^{s \wedge c_0} \cdot \epsilon^{(s - c_0) \vee 0} \\ & \leq \left(\frac{4e(s + s^* - 1)}{s} \right)^{s/2} \cdot \sqrt{s + s^* - 1} \cdot 4 \|\psi\|_2^2 \cdot \epsilon_1^{(s^* - 1 - \lfloor s/2 \rfloor) \vee 0} \cdot \epsilon_0^{s \wedge c_0} \cdot \epsilon^{(s - c_0) \vee 0} \\ & \leq (4es^*)^{s/2} \cdot \sqrt{s + s^*} \cdot 4 \|\psi\|_2^2 \cdot \epsilon_1^{(s^* - 1 - \lfloor s/2 \rfloor) \vee 0} \cdot \epsilon_0^{s \wedge c_0} \cdot \epsilon^{(s - c_0) \vee 0}. \end{aligned}$$

1693 **Here, the second inequality follows from the Stirling's approximation, and the last inequality holds**
 1694 **because $s \geq 1$.**

1695 **Case $s = 0$.** For the case $s = 0$, we have that

$$|\mathbb{E}_{z \sim \mathcal{N}_d} [\psi(\langle w, z \rangle) \psi(\langle w', z \rangle)]| = \left| C(s^*, 0, 0, \langle w, w' \rangle) \cdot \langle w, w' \rangle^{s^* - 1} \right| \leq 4 \|\psi\|_2^2 \epsilon_1^{s^* - 1},$$

1696 **which can also be upper bounded by the quantity derived for the case $s \geq 1$. Hence, the proof is**
 1697 **completed. \square**

1698 **Proposition I.4.** *Let $\psi(\cdot, \cdot)^2$ satisfies the quadratic integrability condition and high-pass condition*
 1699 *in Assumption 4.1 with s^* being the generative exponent. Suppose the remaining definitions ($c_0,$*
 1700 *$\{T_s\}_{s=0}^\infty, \epsilon, \epsilon_0, \epsilon_1$) and conditions are the same as in Lemma I.3. Suppose $c_0 \in \{0, 1, 2\}$ and $s_0 \in \mathbb{N}_0$.*
 1701 *Take function series $\{\zeta_s(\cdot)\}_{s=0}^\infty$ satisfying $\mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] < C, \forall s \in \mathbb{N}_0$ for some universal $C = O(1)$.*
 1702 *Suppose $4e\epsilon^2 < 1/2$. If $s^* = 1$, then we have that*

$$\sum_{s \geq s_0} (s + 2) \cdot |\mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \mathbf{h}_s(z) [T_s]]| \lesssim \epsilon_0^{s_0} + \epsilon_0^{c_0} \epsilon^{(s_0 - c_0) \vee 0}.$$

1703 *If $s^* \geq 2$, then we have that*

$$\begin{aligned} & \sum_{s \geq s_0} (s + 2) \cdot |\mathbb{E}_{\mathbb{Q}} [\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \mathbf{h}_s(z) [T_s]]| \\ & \lesssim \mathbb{1}(s_0 \leq c_0) \cdot \left(\epsilon_1^{s^* - 1 - \lfloor s_0/2 \rfloor} \cdot \epsilon_0^{s_0} + \epsilon_1^{s^* - 1 - \lfloor c_0/2 \rfloor} \cdot \epsilon_0^{c_0} \right) + \epsilon_0^{c_0} \cdot \epsilon^{(2s^*) \vee s_0 - c_0} \\ & \quad + \mathbb{1}(s_0 \leq 2(s^* - 1)) \cdot \left(\epsilon_1^{s^* - 1 - \lfloor (c_0 + 1)/2 \rfloor} \cdot \epsilon_0^{c_0} \cdot \epsilon + \epsilon_0^{c_0} \cdot \epsilon^{2(s^* - 1) + 1 - c_0} \right). \end{aligned}$$

1704 *Here, \lesssim hides constants that depend on $s_0, c_0, C, \mathbb{E}_{\mathbb{Q}}[\psi(y, x)^4]$.*

1705 **Proof of Proposition I.4.** Let F denote the target quantity. Invoking Lemma I.3 for each degree s , we
 1706 **have that**

$$\begin{aligned} F & \leq \sum_{s \geq s_0} \sqrt{s(s+1)} \cdot \mathbb{E}_{y \sim \mathbb{Q}} |\mathbb{E}_{z \sim \mathcal{N}_d} [\zeta_s(y) \psi(y, \langle w, z \rangle) \psi(y, \langle w', z \rangle) \mathbf{h}_s(z) [T_s]]| \\ & \leq \sum_{s \geq s_0} \sqrt{s(s+1)} \cdot \mathbb{E}_{y \sim \mathbb{Q}} [4 |\zeta_s(y)| \mathbb{E}_{x \sim \mathcal{N}} [\psi(y, x)^2]] (4es^*)^{s/2} \sqrt{s + s^*} \epsilon_1^{(s^* - 1 - \lfloor s/2 \rfloor) \vee 0} \epsilon_0^{s \wedge c_0} \epsilon^{(s - c_0) \vee 0} \\ & \leq \sum_{s \geq s_0} \sqrt{s(s+1)} \cdot 4 \sqrt{\mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \mathbb{E}_{\mathbb{Q}}[\psi(y, x)^4]} (4es^*)^{s/2} \sqrt{s + s^*} \epsilon_1^{(s^* - 1 - \lfloor s/2 \rfloor) \vee 0} \epsilon_0^{s \wedge c_0} \epsilon^{(s - c_0) \vee 0}, \end{aligned}$$

1707 **where the last inequality follows from the Cauchy-Schwarz inequality. Noting that by our assumptions,**
 1708 **$\mathbb{E}_{\mathbb{Q}}[\zeta_s(y)^2] \mathbb{E}_{\mathbb{Q}}[\psi(y, x)^4] = O(1)$ uniformly over s , which gives us**

$$F \lesssim \sum_{s \geq s_0} \sqrt{s(s+1)(s+s^*)} (4es^*)^{s/2} \epsilon_1^{(s^* - 1 - \lfloor s/2 \rfloor) \vee 0} \epsilon_0^{s \wedge c_0} \epsilon^{(s - c_0) \vee 0}.$$

1709 For the case $s^* = 1$, we have the above term controlled by

$$\begin{aligned}
F|_{s^*=1} &\lesssim \sum_{s \geq s_0} \sqrt{s(s+1)^2} (4e)^{s/2} \epsilon_0^{s \wedge c_0} \epsilon^{(s-c_0) \vee 0} \\
&\lesssim \mathbf{1}(c_0 > s_0) \cdot \sup_{s_0 \leq s \leq c_0-1} \epsilon_0^s + \sum_{s \geq c_0 \vee s_0} \sqrt{s(s+1)^2} (4e)^{s/2} \epsilon_0^{c_0} \epsilon^{s-c_0} \\
&\lesssim \mathbf{1}(c_0 > s_0) \cdot \epsilon_0^{s_0} + \epsilon_0^{c_0} \epsilon^{(s_0-c_0) \vee 0} \sum_{s \geq c_0 \vee s_0} \sqrt{s(s+1)^2} (4e\epsilon^2)^{(s-c_0 \vee s_0)/2} \\
&\lesssim \mathbf{1}(s_0 < c_0) \cdot \epsilon_0^{s_0} + \epsilon_0^{c_0} \epsilon^{(s_0-c_0) \vee 0} \lesssim \epsilon_0^{c_0 \wedge s_0} \epsilon^{(s_0-c_0) \vee 0},
\end{aligned}$$

1710 where \lesssim only hides constants that depend on s_0, c_0 . The last second inequality holds by noting that
1711 $4e\epsilon^2 < 1/2$.

1712 For the case $s^* \geq 2$, we note that $c_0 \leq 2 \leq 2(s^* - 1)$, and we have

$$\begin{aligned}
F|_{s^* \geq 2} &\lesssim \sum_{s \geq s_0} \sqrt{s(s+1)(s+s^*)} (4es^*)^{s/2} \epsilon_1^{(s^*-1-\lfloor s/2 \rfloor) \vee 0} \epsilon_0^{s \wedge c_0} \epsilon^{(s-c_0) \vee 0} \\
&\lesssim \mathbf{1}(s_0 \leq c_0) \cdot \max_{s_0 \leq s \leq c_0} \epsilon_1^{s^*-1-\lfloor s/2 \rfloor} \cdot \epsilon_0^s \\
&\quad + \mathbf{1}(s_0 \leq 2(s^* - 1)) \cdot \max_{c_0+1 \leq s \leq 2(s^*-1)+1} \epsilon_1^{s^*-1-\lfloor s/2 \rfloor} \cdot \epsilon_0^{c_0} \cdot \epsilon^{s-c_0} \\
&\quad + \sum_{s \geq (2(s^*-1)+2) \vee s_0} \sqrt{s(s+1)(s+s^*)} (4es^*)^{s/2} \epsilon_0^{c_0} \epsilon^{s-c_0} \\
&\lesssim \mathbf{1}(s_0 \leq c_0) \cdot \left(\epsilon_1^{s^*-1-\lfloor s_0/2 \rfloor} \cdot \epsilon_0^{s_0} + \epsilon_1^{s^*-1-\lfloor c_0/2 \rfloor} \cdot \epsilon_0^{c_0} \right) + \epsilon_0^{c_0} \cdot \epsilon^{(2(s^*-1)+2) \vee s_0 - c_0} \\
&\quad + \mathbf{1}(s_0 \leq 2(s^* - 1)) \cdot \left(\epsilon_1^{s^*-1-\lfloor (c_0+1)/2 \rfloor} \cdot \epsilon_0^{c_0} \cdot \epsilon + \epsilon_0^{c_0} \cdot \epsilon^{2(s^*-1)+1-c_0} \right).
\end{aligned}$$

1713 Thus, we complete the proof. \square

1714 I.2 Technical Results for Uniform Distribution on the Sphere

1715 **Lemma I.5** (Moment of polynomial on a sphere, adapted from Folland (2001)). *Let $\xi =$
1716 $(\xi_1, \xi_2, \dots, \xi_d) \sim \text{Unif}(\mathbb{S}^{d-1})$ and $s_1, s_2, \dots, s_d \in \mathbb{N}_0$. Let $s = \sum_{i=1}^d s_i$. Then we have*

$$\mathbb{E}_\xi \left[\prod_{i=1}^d \xi_i^{s_i} \right] = \begin{cases} 0 & \text{if some } s_i \text{ is odd,} \\ \prod_{i=1}^d \frac{\Gamma((s_i+1)/2)}{\Gamma(1/2)} \cdot \frac{\Gamma(d/2)}{\Gamma((s+d)/2)} & \text{if all } s_i \text{ are even,} \end{cases}$$

1717 To evaluate this moment, we provide the following bound.

1718 **Fact I.6.** *Take even degrees $s_1, s_2, \dots, s_d \in \mathbb{N}_0$ and $s = \sum_{i=1}^d s_i$. Then we have*

$$\left(\frac{1}{d} \right)^{s/2} \leq \prod_{i=1}^d \frac{\Gamma((s_i+1)/2)}{\Gamma(1/2)} \cdot \frac{\Gamma(d/2)}{\Gamma((s+d)/2)} \leq \left(\frac{s}{d} \right)^{s/2},$$

1719 where we define $(0)^0 = 1$.

1720 *Proof of Fact I.6.* Note that we can rewrite the product as

$$\prod_{i=1}^d \frac{\Gamma((s_i+1)/2)}{\Gamma(1/2)} \cdot \frac{\Gamma(d/2)}{\Gamma((s+d)/2)} = \frac{\prod_{i=1}^d (s_i-1)!!}{(s+d-2) \cdot (s+d-4) \cdots d}.$$

1721 Using the fact that $a/b \geq (a-1)/(b-1)$ for $2 \leq a \leq b$, we can recursively apply this inequality to
1722 the factorial until it gives us the desired lower bound of $(1/d)^{s/2}$. For the upper bound, we can lower
1723 bound the denominator by $d^{s/2}$ and upper bound the numerator by $s^{s/2}$. We thus have the desired
1724 result. \square

1725 Note that the second moment $\mathbb{E}[\xi_i^2] \asymp d^{-1}$, thus Fact I.6 can be viewed as some kind of *reverse*
 1726 *Holder's inequality* where we use lower moment to control higher moment. Notably, the reverse
 1727 inequality gives a *dimension-free* bound for the moments of the polynomial on the sphere. The
 1728 following proposition formalizes this intuition.

1729 **Proposition I.7.** *Suppose $\xi = (\xi_1, \xi_2, \dots, \xi_d) \sim \text{Unif}(\mathbb{S}^{d-1})$. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function such*
 1730 *that $\mathbb{E}[(f(\xi) - 1)^2] \leq \varepsilon^2$. Take nonnegative degree $\mathbf{s} = (s_1, s_2, \dots, s_d)$ with $\|\mathbf{s}\|_1 = s$, and suppose*
 1731 *each s_i is even for $i \in [d]$. We then have*

$$(1 - (2s)^{s/2}\varepsilon) \cdot \left(\frac{1}{d}\right)^{s/2} \leq \mathbb{E} \left[f(\xi) \cdot \prod_{i=1}^d \xi_i^{s_i} \right] \leq (1 + 2^{s/2}\varepsilon) \cdot \left(\frac{s}{d}\right)^{s/2}.$$

1732 *Proof of Proposition I.7.* By the Cauchy-Schwarz inequality, we have for any $j = 0, 1, \dots, s_1/2$,

$$\left| \mathbb{E} \left[(f(\xi) - 1) \cdot \prod_{i=1}^d \xi_i^{s_i} \right] \right| \leq \sqrt{\mathbb{E}[(f(\xi) - 1)^2]} \cdot \sqrt{\mathbb{E} \left[\prod_{i=1}^d \xi_i^{2s_i} \right]} \leq \varepsilon \cdot \left(\frac{2s}{d}\right)^{s/2},$$

1733 where we use Lemma I.5 and the upper bound in Fact I.6 for the second inequality. Additionally, note
 1734 that each s_i is even, we use the same argument to have

$$\left(\frac{1}{d}\right)^{s/2} \leq \mathbb{E} \left[\prod_{i=1}^d \xi_i^{s_i} \right] \leq \left(\frac{s}{d}\right)^{s/2}.$$

1735 Combining these two inequalities, we conclude the proof. \square

1736 **Proposition I.8.** *Let $\xi = (\xi_1, \xi_2, \dots, \xi_d) \sim \text{Unif}(\mathbb{S}^{d-1})$, $\gamma > 0$ be a fixed parameter, and $\mathbf{s} =$
 1737 (s_1, s_2, \dots, s_d) be a nonnegative integer vector with $\|\mathbf{s}\|_1 = s$. It then holds that*

$$\begin{aligned} G(\mathbf{s}) &:= \mathbb{E} \left[(\xi_1 + \gamma)^{s_1} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] = \mathbb{E} \left[\sum_{j=0}^{\lfloor s_1/2 \rfloor} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \xi_1^{2j} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \\ &= \begin{cases} 0, & \text{if some } s_i \text{ is odd for } i = 2, 3, \dots, d, \\ C(\mathbf{s}, \gamma) \cdot \left(\gamma + \frac{1}{\sqrt{d}}\right)^{2\lfloor s_1/2 \rfloor} \cdot \gamma^{\mathbf{1}(s_1 \text{ odd})} \cdot \left(\frac{1}{\sqrt{d}}\right)^{s-s_1}, & \text{otherwise,} \end{cases} \end{aligned}$$

1738 where $1/5 \leq C(\mathbf{s}, \gamma) \leq s^{s/2}$.

1739 *Proof of Proposition I.8.* Note that

$$G(\mathbf{s}) = \mathbb{E} \left[(\xi_1 + \gamma)^{s_1} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] = \mathbb{E} \left[\sum_{j=0}^{s_1} \binom{s_1}{j} \xi_1^j \gamma^{s_1-j} \cdot \prod_{i=2}^d \xi_i^{s_i} \right].$$

1740 Note that if there exists any odd degree s_i for $i \geq 2$, then $G(\mathbf{s}) = 0$. For s_2, s_3, \dots, s_d being even,
 1741 we have

$$\begin{aligned} G(\mathbf{s}) &= \mathbb{E} \left[\sum_{j=0}^{\lfloor s_1/2 \rfloor} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \xi_1^{2j} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \\ &= \sum_{j=0}^{\lfloor s_1/2 \rfloor} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \frac{\Gamma(j+1/2)}{\Gamma(1/2)} \cdot \prod_{i=2}^d \frac{\Gamma((s_i+1)/2)}{\Gamma(1/2)} \cdot \frac{\Gamma(d/2)}{\Gamma((s-s_1+2j+d)/2)} \\ &= \sum_{j=0}^{\lfloor s_1/2 \rfloor} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \underbrace{\frac{(2j-1)!! \cdot \prod_{i=2}^d (s_i-1)!!}{(s-s_1+2j+d-2) \cdot (s-s_1+2j+d-4) \cdots d}}_{(I)}, \end{aligned}$$

1742 where for the second identity, we use Lemma I.5 to compute the moments of the polynomial on the
 1743 sphere. Now, we compute the lower and upper bounds of $G(\mathbf{s})$.

1744 **Lower Bound.** For the lower bound, we have

$$\begin{aligned}
G(\mathbf{s}) &\geq \sum_{j=0}^{\lfloor s_1/2 \rfloor} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot d^{-j-(s_2+s_3+\dots+s_d)/2} \\
&\geq \frac{1}{5} \cdot 5 \sum_{j=0}^{\lfloor s_1/2 \rfloor} \binom{2\lfloor s_1/2 \rfloor}{2j} \cdot \gamma^{2\lfloor s_1/2 \rfloor-2j} \left(\frac{1}{\sqrt{d}}\right)^{2j} \cdot d^{-(s_2+s_3+\dots+s_d)/2} \cdot \gamma^{\mathbb{1}(s_1 \text{ odd})} \\
&\geq \frac{1}{5} \cdot \sum_{j=0}^{2\lfloor s_1/2 \rfloor} \binom{2\lfloor s_1/2 \rfloor}{j} \cdot \gamma^{2\lfloor s_1/2 \rfloor-j} \left(\frac{1}{\sqrt{d}}\right)^j \cdot d^{-(s_2+s_3+\dots+s_d)/2} \cdot \gamma^{\mathbb{1}(s_1 \text{ odd})}.
\end{aligned}$$

1745 Here in the first inequality, we use the fact that $a/b \geq (a-1)/(b-1)$ for $2 \leq a \leq b$ and apply it
1746 recursively to the factorial (I) until it gives us $d^{-j-(s_2+s_3+\dots+s_d)/2}$. For the second inequality, we
1747 first rearrange the terms in the summation and lower bound the binomial coefficients by changing s_1
1748 to $2\lfloor s_1/2 \rfloor$. For the last inequality, let us define

$$A_j = \binom{2\lfloor s_1/2 \rfloor}{j} \cdot \gamma^{2\lfloor s_1/2 \rfloor-j} \left(\frac{1}{\sqrt{d}}\right)^j.$$

1749 We invoke Lemma J.4 and have that for each odd j , we have that

$$A_j \leq 2(A_{j-1} + A_{j+1}), \quad j = 1, 3, \dots, 2\lfloor s_1/2 \rfloor - 1.$$

1750 Therefore, the summation of all the odd terms is upper bounded by 4 times the summation of all the
1751 even terms, which gives us the last inequality. Therefore, the lower bound of $G(\mathbf{s})$ is

$$G(\mathbf{s}) \geq \frac{1}{5} \cdot \left(\gamma + \frac{1}{\sqrt{d}}\right)^{2\lfloor s_1/2 \rfloor} \cdot \gamma^{\mathbb{1}(s_1 \text{ odd})} \cdot \left(\frac{1}{\sqrt{d}}\right)^{(s-s_1)}.$$

1752 **Upper Bound.** For the upper bound, we have

$$\begin{aligned}
G(\mathbf{s}) &\leq \sum_{j=0}^{\lfloor s_1/2 \rfloor} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot d^{-j-(s_2+s_3+\dots+s_d)/2} \cdot s^{s/2} \\
&\leq \sum_{j=0}^{2\lfloor s_1/2 \rfloor} \binom{2\lfloor s_1/2 \rfloor}{j} \cdot \gamma^{2\lfloor s_1/2 \rfloor-j} \left(\frac{1}{\sqrt{d}}\right)^j \cdot d^{-(s_2+s_3+\dots+s_d)/2} \cdot \gamma^{\mathbb{1}(s_1 \text{ odd})} \cdot s^{s/2} \\
&= s^{s/2} \cdot \left(\gamma + \frac{1}{\sqrt{d}}\right)^{2\lfloor s_1/2 \rfloor} \cdot \gamma^{\mathbb{1}(s_1 \text{ odd})} \cdot \left(\frac{1}{\sqrt{d}}\right)^{-(s-s_1)},
\end{aligned}$$

1753 where in the first line, we lower bound the denominator in the factorial (I) by $d^{j+(s_2+s_3+\dots+s_d)/2}$
1754 and upper bound the numerator by $s^{s/2}$. For the second inequality, we use the nonnegativity of each
1755 terms and append the terms with j being odd to the summation.

1756 Combining the lower and upper bounds, we have that

$$G(\mathbf{s}) = \begin{cases} 0, & \text{if some } s_i \text{ is odd for } i = 2, 3, \dots, d, \\ C(\mathbf{s}, \gamma) \cdot \left(\gamma + \frac{1}{\sqrt{d}}\right)^{2\lfloor s_1/2 \rfloor} \cdot \gamma^{\mathbb{1}(s_1 \text{ odd})} \cdot \left(\frac{1}{\sqrt{d}}\right)^{-(s-s_1)}, & \text{otherwise,} \end{cases}$$

1757 for $1/5 \leq C(\mathbf{s}, \gamma) \leq s^{s/2}$. Hence, the proof is complete. \square

1758 I.3 Technical Results on Polarized Random Vectors

1759 **Lemma I.9** (Moments of weakly polarized random vector). *Suppose $\xi = (\xi_1, \xi_2, \dots, \xi_d) \sim$
1760 $\text{Unif}(\mathbb{S}^{d-1})$ and define the polarized vector w as*

$$w = \frac{\xi + \gamma e_1}{\|\xi + \gamma e_1\|_2},$$

1761 where $\gamma = o(1) > 0$ is a parameter that describes the polarization strength, and $e_1 = (1, 0, \dots, 0)$
 1762 is the first standard basis vector. Take nonnegative integer degree $\mathbf{s} = (s_1, s_2, \dots, s_d)$ with $\|\mathbf{s}\|_1 =$
 1763 $s = O(1) < (2\sqrt{e\gamma})^{-1} - 2$. Then we have

$$\mathbb{E}_w \left[\prod_{i=1}^d w_i^{s_i} \right] = \begin{cases} 0, & \text{if some } s_i \text{ is odd for } i = 2, 3, \dots, d, \\ C(\mathbf{s}, \gamma) \cdot \left(\gamma + \frac{1}{\sqrt{d}} \right)^{s_1} \cdot \left(\frac{1}{\sqrt{d}} \right)^{s-s_1}, & \text{if } s_1 \text{ is even, } s_2, \dots, s_d \text{ are even,} \\ C(\mathbf{s}, \gamma) \cdot \gamma \cdot \left(\gamma + \frac{1}{\sqrt{d}} \right)^{s_1-1} \cdot \left(\frac{1}{\sqrt{d}} \right)^{s-s_1}, & \text{if } s_1 \text{ is odd, } s_2, \dots, s_d \text{ are even,} \end{cases}$$

1764 where $1/5 - O(\gamma) \leq C(\mathbf{s}, \gamma) \leq s^{(s+1)/2} + O(\gamma)$.

1765 *Proof of Lemma I.9.* Let $r = \|\xi + \gamma e_1\|_2$. We reformulate the moment as

$$\mathbb{E} \left[\prod_{i=1}^d w_i^{s_i} \right] = \mathbb{E} \left[\frac{(\xi_1 + \gamma)^{s_1}}{r^s} \cdot \prod_{i=2}^d \xi_i^{s_i} \right].$$

1766 By symmetry, we have the moment equal zero if some s_i is odd for $i = 2, 3, \dots, d$. Hence, we only
 1767 need to consider s_2, s_3, \dots, s_d being even. In the following, we study the case for s_1 being even and
 1768 odd separately.

1769 **Case 1: s_1 is even.** We first look at the simpler case where s_1 is even. We note by weak polarization
 1770 in the sense of $\|\xi\|_2 \gg \|\gamma e_1\|_2$, we can approximate $1/r^s$ as 1 with approximation error

$$\begin{aligned} \left| \mathbb{E} \left[\prod_{i=1}^d w_i^{s_i} \right] - \mathbb{E} \left[(\xi_1 + \gamma)^{s_1} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \right| &\leq \sqrt{\mathbb{E}[(1-r^{-s})^2]} \cdot \sqrt{\mathbb{E} \left[(\xi_1 + \gamma)^{2s_1} \cdot \prod_{i=2}^d \xi_i^{2s_i} \right]} \\ &\leq es\gamma \cdot \sqrt{\mathbb{E} \left[(\xi_1 + \gamma)^{2s_1} \cdot \prod_{i=2}^d \xi_i^{2s_i} \right]}, \end{aligned} \quad (I.4)$$

1771 where we use the Cauchy-Schwarz inequality in the first line and for the second line, we use the fact

$$|1 - r^{-s}| \leq \frac{1}{(1-\gamma)^s} - 1 = \left(1 + \frac{\gamma}{1-\gamma} \right)^s - 1 \leq \exp\left(\frac{s\gamma}{1-\gamma}\right) - 1 \leq es\gamma \quad (I.5)$$

1772 for $s < (2\sqrt{e\gamma})^{-1}$. Define

$$G(\mathbf{s}) := \mathbb{E} \left[(\xi_1 + \gamma)^{s_1} \cdot \prod_{i=2}^d \xi_i^{s_i} \right].$$

1773 By Proposition I.8, we have that for even s_1, s_2, \dots, s_d ,

$$G(\mathbf{s}) = C'(\mathbf{s}, \gamma) \cdot \left(\gamma + \frac{1}{\sqrt{d}} \right)^{s_1} \cdot \left(\frac{1}{\sqrt{d}} \right)^{s-s_1},$$

1774 with $1/5 \leq C'(\mathbf{s}, \gamma) \leq s^{s/2}$. Plugging the form of $G(\mathbf{s})$ into (I.4), we have that

$$\begin{aligned} \mathbb{E}_w \left[\prod_{i=1}^d w_i^{s_i} \right] &= M(\mathbf{s}) \pm es\gamma \cdot \sqrt{M(2\mathbf{s})} \\ &= \left(C'(\mathbf{s}, \gamma) \pm es\gamma \sqrt{C'(2\mathbf{s}, \gamma)} \right) \cdot \left(\gamma + \frac{1}{\sqrt{d}} \right)^{s_1} \cdot \left(\frac{1}{\sqrt{d}} \right)^{s-s_1} \\ &= \left(C'(\mathbf{s}, \gamma) \pm es\gamma (2s)^{s/2} \right) \cdot \left(\gamma + \frac{1}{\sqrt{d}} \right)^{s_1} \cdot \left(\frac{1}{\sqrt{d}} \right)^{s-s_1}. \end{aligned}$$

1775 Here, $C'(\mathbf{s}, \gamma) - es\gamma (2s)^{s/2} \geq 1/5 - es\gamma (2s)^{s/2}$ and $C'(\mathbf{s}, \gamma) + es\gamma (2s)^{s/2} \leq s^{s/2} + es\gamma (2s)^{s/2}$.

1776 **Case 2: s_1 is odd.** We now consider the more complicated case where s_1 is odd. In the following
 1777 proof, we will frequently invoke the following proposition, whose proof is deferred to Proposition I.7
 1778 of Appendix H.

1779 **Proposition I.7 (Restated).** Suppose $\xi = (\xi_1, \xi_2, \dots, \xi_d) \sim \text{Unif}(\mathbb{S}^{d-1})$. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a
 1780 function such that $\mathbb{E}[(f(\xi) - 1)^2] \leq \varepsilon^2$. Take nonnegative degree $\mathbf{s} = (s_1, s_2, \dots, s_d)$ with $\|\mathbf{s}\|_1 = s$,
 1781 and suppose each s_i is even for $i \in [d]$. We then have

$$(1 - (2s)^{s/2}\varepsilon) \cdot \left(\frac{1}{d}\right)^{s/2} \leq \mathbb{E} \left[f(\xi) \cdot \prod_{i=1}^d \xi_i^{s_i} \right] \leq (1 + 2^{s/2}\varepsilon) \cdot \left(\frac{s}{d}\right)^{s/2}.$$

1782 We first rewrite the moment as

$$\begin{aligned} \mathbb{E} \left[\prod_{i=1}^d w_i^{s_i} \right] &= \mathbb{E} \left[\frac{(\xi_1 + \gamma)^{s_1}}{r^s} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] & (I.6) \\ &= \underbrace{\sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \mathbb{E} \left[\frac{\xi_1^{2j}}{r^s} \cdot \prod_{i=2}^d \xi_i^{s_i} \right]}_{(I) \text{ even terms}} + \underbrace{\sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j+1} \gamma^{s_1-2j-1} \cdot \mathbb{E} \left[\frac{\xi_1^{2j+1}}{r^s} \cdot \prod_{i=2}^d \xi_i^{s_i} \right]}_{(II) \text{ odd terms}}. \end{aligned}$$

1783 Let us look at the odd terms of (I.6). Let $r_+ = \|\xi + \gamma e_1\|_2$ and $r_- = \|-\xi + \gamma e_1\|_2$. By symmetry,
 1784 we have for the expectation within the odd terms that

$$\begin{aligned} \mathbb{E} \left[\frac{\xi_1^{2j+1}}{r^s} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] &= \frac{1}{2} \cdot \mathbb{E} \left[\left(\frac{1}{r_+^s} - \frac{1}{r_-^s} \right) \cdot \xi_1^{2j+1} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \\ &= \frac{1}{2} \cdot \mathbb{E} \left[\frac{(r_-^2 - r_+^2) \left(\sum_{l=0}^{s-1} r_-^l r_+^{s-1-l} \right)}{r_+^s \cdot r_-^s \cdot (r_+ + r_-)} \cdot \xi_1^{2j+1} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \\ &= -\gamma \cdot \mathbb{E} \left[\frac{2 \cdot \sum_{l=0}^{s-1} r_-^l r_+^{s-1-l}}{r_+^s \cdot r_-^s \cdot (r_+ + r_-)} \cdot \xi_1^{2j+2} \cdot \prod_{i=2}^d \xi_i^{s_i} \right]. \end{aligned}$$

1785 Here, the last identity holds by noting that $r_-^2 - r_+^2 = -4\gamma\xi_1$. Note that

$$\sup_{\xi \in \mathbb{S}^{d-1}} \left| \frac{2 \cdot \sum_{l=0}^{s-1} r_-^l r_+^{s-1-l}}{r_+^s \cdot r_-^s \cdot (r_+ + r_-)} - 1 \right| \leq \left(\frac{1}{1-\gamma} \right)^{s+2} - 1 \leq e\gamma(s+2).$$

1786 Hence, we have for the odd terms (II) that

$$\begin{aligned} -(II) &= \mathbb{E} \left[\frac{2 \cdot \sum_{l=0}^{s-1} r_-^l r_+^{s-1-l}}{r_+^s \cdot r_-^s \cdot (r_+ + r_-)} \cdot \sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j+1} \gamma^{s_1-2j} \cdot \xi_1^{2j+2} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \\ &\leq \sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j+1} \gamma^{s_1-2j} \cdot (1 + 2^{(s+2)/2} e\gamma(s+2)) \cdot \left(\frac{s+2}{d} \right)^{(s-s_1+2j+2)/2} \\ &\leq s_1 \cdot (1 + 2^{(s+2)/2} e\gamma(s+2)) \cdot \left(\frac{s+2}{d} \right)^{(s-s_1+2)/2} \cdot \gamma \cdot \sum_{j=0}^{s_1-1} \binom{s_1-1}{j} \cdot \gamma^{s_1-1-j} \cdot \left(\sqrt{\frac{s+2}{d}} \right)^j, \end{aligned}$$

1787 where the first inequality holds by Proposition I.7 and the second inequality holds by $\binom{s_1}{2j+1} / \binom{s_1-1}{2j} =$
 1788 $s_1 / (2j+1) \leq s_1$ and also appending the odd terms to the summation. Thus, we have

$$-(II) \leq s_1 \cdot (1 + 2^{(s+2)/2} e\gamma(s+2)) \cdot \gamma \cdot \left(\gamma + \sqrt{\frac{s+2}{d}} \right)^{s_1-1} \cdot \left(\sqrt{\frac{s+2}{d}} \right)^{s-s_1+2}.$$

1789 Next, we study the even terms (I) in (I.6). Using Proposition I.7 with the uniform bound in (I.5), we
 1790 have (I) upper bounded by

$$\begin{aligned}
 \text{(I)} &= \sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \mathbb{E} \left[\frac{\xi_1^{2j}}{r^{2j}} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \\
 &\leq (1 + 2^{s/2} e s \gamma) \cdot \sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \left(\sqrt{\frac{s}{d}} \right)^{s-s_1+2j} \\
 &\leq (1 + 2^{s/2} e s \gamma) \cdot s_1 \cdot \sum_{j=0}^{s_1-1} \binom{s_1-1}{j} \cdot \gamma^{s_1-1-j} \cdot \left(\sqrt{\frac{s}{d}} \right)^j \cdot \left(\sqrt{\frac{s}{d}} \right)^{s-s_1} \cdot \gamma \\
 &= s_1 \cdot (1 + 2^{s/2} e s \gamma) \cdot \gamma \cdot \left(\gamma + \sqrt{\frac{s}{d}} \right)^{s_1-1} \cdot \left(\sqrt{\frac{s}{d}} \right)^{s-s_1}.
 \end{aligned}$$

1791 Similarly, we have the lower bound for (I) as

$$\begin{aligned}
 \text{(I)} &= \sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \mathbb{E} \left[\frac{\xi_1^{2j}}{r^{2j}} \cdot \prod_{i=2}^d \xi_i^{s_i} \right] \\
 &\geq (1 - (2s)^{s/2} e s \gamma) \cdot \sum_{j=0}^{(s_1-1)/2} \binom{s_1}{2j} \gamma^{s_1-2j} \cdot \left(\sqrt{\frac{1}{d}} \right)^{s-s_1+2j} \\
 &\geq \frac{1}{5} (1 - (2s)^{s/2} e s \gamma) \cdot \sum_{j=0}^{s_1-1} \binom{s_1-1}{j} \gamma^{s_1-1-j} \cdot \left(\sqrt{\frac{1}{d}} \right)^j \cdot \left(\sqrt{\frac{1}{d}} \right)^{s-s_1} \cdot \gamma \\
 &= \frac{1}{5} (1 - (2s)^{s/2} e s \gamma) \cdot \gamma \cdot \left(\gamma + \sqrt{\frac{1}{d}} \right)^{s_1-1} \cdot \left(\sqrt{\frac{1}{d}} \right)^{s-s_1}.
 \end{aligned}$$

1792 Combining these upper and lower bounds for (I) together with the upper bound for (II), we have

$$\begin{aligned}
 \mathbb{E} \left[\prod_{i=1}^d w_i^{s_i} \right] &\geq \left(\frac{1}{5} (1 - (2s)^{s/2} e s \gamma) - \frac{s_1 \cdot (1 + 2^{(s+2)/2} e \gamma (s+2)) \cdot (s+2)^{(s+2)/2}}{d} \right) \\
 &\quad \cdot \gamma \cdot \left(\gamma + \sqrt{\frac{1}{d}} \right)^{s_1-1} \cdot \left(\sqrt{\frac{1}{d}} \right)^{s-s_1},
 \end{aligned}$$

1793 and

$$\mathbb{E} \left[\prod_{i=1}^d w_i^{s_i} \right] \leq (1 + 2^{s/2} e s \gamma) \cdot s^{(s+1)/2} \cdot \gamma \cdot \left(\gamma + \sqrt{\frac{1}{d}} \right)^{s_1-1} \cdot \left(\sqrt{\frac{1}{d}} \right)^{s-s_1}.$$

1794 Hence, we complete our proof. \square

1795 J Auxiliary Lemmas

1796 **Lemma J.1** (Bernstein's inequality). *Let X_1, \dots, X_n be independent random variables with $|X_i -$
 1797 $\mathbb{E}[X_i]| \leq C$ for all $i \in [n]$. Then for any $t > 0$, it holds that*

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E} X_i \right| \geq t \right) \leq 2 \exp \left(- \frac{nt^2/2}{n^{-1} \cdot \sum_{i=1}^n \text{Var}[X_i] + Ct/3} \right),$$

1798 or equivalently, for any $\delta \in (0, 1)$,

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E} X_i \right| \leq \sqrt{\frac{2 \cdot n^{-1} \sum_{i=1}^n \text{Var}[X_i] \cdot \log \delta^{-1}}{n} + \frac{C \log \delta^{-1}}{3n}} \right) \geq 1 - \delta.$$

1799 For the vector case, by a union bound over all the coordinates, we have the following corollary.

1800 **Corollary J.2** (Vector version of Bernstein's inequality). *Let X_1, \dots, X_n be independent random*
 1801 *vectors in \mathbb{R}^d with $\|X_i - \mathbb{E}[X_i]\|_\infty \leq C$ for all $i \in [n]$. Then for any $\delta \in (0, 1)$, it holds with*
 1802 *probability at least $1 - \delta$ that*

$$\left\| \frac{1}{n} \sum_{i=1}^n X_i \right\|_2 \lesssim \left\| \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n X_i \right] \right\|_2 + \sqrt{\frac{n^{-1} \sum_{i=1}^n \text{Tr}(\text{Cov}[X_i]) \cdot \log(d\delta^{-1})}{n}} + \frac{\sqrt{d}C \log(d\delta^{-1})}{n}.$$

1803 **Lemma J.3** (Lemma I.3. in Damian et al. (2024)). *Let $X_1, \dots, X_n \in \mathbb{R}^d$ be independent mean-zero*
 1804 *random vectors such that for all $p \geq 2$, $\mathbb{E}[\|X_i\|^p]^{1/p} \leq Cp^{k/2}$ for some constants $k, C \geq 0$ and*
 1805 *vector norm $\|\cdot\|$. Define $\sigma^2 := n^{-1} \sum_{i=1}^n \mathbb{E}[\|X_i\|^2]$ and $Y := n^{-1} \cdot \sum_{i=1}^n X_i$. Then with probability*
 1806 *at least $1 - 2\delta$,*

$$\|Y\| \lesssim \sigma \cdot \sqrt{\frac{\log(1/\delta)}{n}} + \frac{C \log(1/\delta) \log(n/\delta)^{k/2}}{n}.$$

1807 where \lesssim only hides constant that depends on k .

1808 **Lemma J.4** (Ratio bound of binomial expansion). *For real numbers $a, b > 0$ and integer $s \geq 2$,*
 1809 *define*

$$A_j := \binom{s}{j} \cdot a^{s-j} \cdot b^j, \quad \text{for } j = 0, 1, \dots, s.$$

1810 Then for all $j = 1, 2, \dots, s-1$, it holds that $A_j \leq 2(A_{j-1} + A_{j+1})$.

1811 *Proof of Lemma J.4.* By the definition of A_j ,

$$\begin{aligned} \min \left\{ \frac{A_j}{A_{j-1}}, \frac{A_j}{A_{j+1}} \right\} &= \min \left\{ \frac{b}{a} \cdot \frac{s-j+1}{j}, \frac{a}{b} \cdot \frac{j+1}{s-j} \right\} \leq \sqrt{\frac{b}{a} \cdot \frac{s-j+1}{j} \cdot \frac{a}{b} \cdot \frac{j+1}{s-j}} \\ &\leq \sqrt{\left(1 + \frac{1}{s-j}\right) \left(1 + \frac{1}{j}\right)} \leq 2. \end{aligned}$$

1812 Hence, the proof is complete by the nonnegativity of A_j . □

1813 **Proposition J.5.** *For $\epsilon \in [0, 1/2)$ and $s, r \in \mathbb{N}_0$ with $s \geq 2r-1$, it holds that*

$$\sum_{j=0}^{\infty} (j+s)(j+s-1) \cdots (j+s-r+1) \cdot \epsilon^j \leq 2s \cdot (s-1) \cdots (s-r+1) \cdot \frac{1}{1-\epsilon}.$$

1814 *Proof of Proposition J.5.* Denote $F(x) = \sum_{j=0}^{\infty} (j+s)(j+s-1) \cdots (j+s-r+1)x^j$ for $x \in (0, 1)$.
 1815 The desired quantity on the left-hand side of the inequality is simply $F(\epsilon)$. It can be verified using
 1816 the expansion of $1/(1-x) = \sum_{j=0}^{\infty} x^j$ for $x \in (0, 1)$ that

$$F(x) = \frac{d^r}{dx^r} \left(\frac{x^s}{1-x} \right) \cdot x^{-(s-r)}.$$

1817 Expanding this expression, we have

$$\begin{aligned} F(x) &= \sum_{\tau=0}^r \binom{r}{\tau} \frac{d^\tau}{dx^\tau} (x^s) \cdot \frac{d^{r-\tau}}{dx^{r-\tau}} \left(\frac{1}{1-x} \right) \cdot x^{-(s-r)} \\ &= \sum_{\tau=0}^r \binom{r}{\tau} s(s-1) \cdots (s-\tau+1) \cdot x^{s-\tau} \cdot \frac{(-1)^{r-\tau} (r-\tau)!}{(1-x)^{r-\tau+1}} \cdot x^{-(s-r)} \\ &= \sum_{\tau=0}^r (r \cdot (r-1) \cdots (\tau+1)) \cdot (s \cdot (s-1) \cdots (s-\tau+1)) \cdot (-1)^{r-\tau} \cdot \frac{x^{r-\tau}}{(1-x)^{r-\tau+1}}. \end{aligned}$$

1818 Write the above summation as $F(x) = \sum_{\tau=0}^r F_{\tau}(x)$, where $F_{\tau}(x)$ is the τ -th term in the summation.
 1819 Then for each $\tau = 0, \dots, r-1$, when $x \in (0, 1)$, we have

$$\frac{F_{\tau}(x)}{F_{\tau+1}(x)} = -\frac{s-\tau}{\tau+1} \cdot \frac{1-x}{x} < -1.$$

1820 Note that $F_r(x)$ is positive, and for each $k = 1, \dots, \lfloor r/2 \rfloor$, $F_{r-2k+1}(x) < -F_{r-2k}(x) < 0$. Since
 1821 $F(x)$ is positive, it is then upper bounded by $F_r(x) + |F_0(x)|$, i.e.,

$$F(x) \leq s \cdot (s-1) \cdots (s-r+1) \cdot \frac{1}{1-x} + r! \cdot \frac{x^r}{(1-x)^{r+1}} \leq 2s \cdot (s-1) \cdots (s-r+1) \cdot \frac{1}{1-x}.$$

1822 The proof is complete by setting $x = \epsilon$. \square

1823 **Lemma J.6** (Gaussian-like tail bound for spherical coordinate). *Suppose $\xi \sim \text{Unif}(\mathbb{S}^{d-1})$. Then the*
 1824 *first coordinate of ξ , denoted by ξ_1 , satisfies that for any $t \geq 0$,*

$$\Pr(\xi_1 \geq \sqrt{Cd^{-1} \log d}) \leq \exp(-d/16) + d^{-C/4},$$

1825 *where $C > 0$ is a constant.*

1826 *Proof of Lemma J.6.* Consider $z \sim \mathcal{N}(0, I_d)$, and it holds that $\xi_1 \stackrel{d}{=} z_1 / \|z\|_2$, where z_1 is the first
 1827 coordinate of z . Note that $\|z\|_2^2 \sim \chi_d^2$, and by a standard tail bound for the χ_d^2 distribution, we have

$$\Pr(\|z\|_2^2 \leq d - 2\sqrt{dx}) \leq \exp(-x), \quad \text{for any } x \geq 0.$$

1828 By taking $x = d/16$, we get $\Pr(\|z\|_2^2 \leq d/2) \leq \exp(-d/16)$. Thus, applying a union bound,

$$\begin{aligned} \Pr(\xi_1 \geq t) &= \Pr\left(\frac{z_1}{\|z\|_2} \geq t\right) \leq \Pr(\|z\|_2^2 \leq d/2) + \Pr(z_1 \geq t\sqrt{d/2}) \\ &\leq \exp(-d/16) + \exp(-t^2 d/4). \end{aligned}$$

1829 The proof is complete by taking $t = \sqrt{Cd^{-1} \log d}$. \square

1830 **Lemma J.7** (Hypergeometric behavior). *Consider random variable $X \sim \text{Hypergeometric}(d, k, k)$*
 1831 *with probability mass*

$$\Pr(X = x) = \frac{\binom{k}{x} \binom{d-k}{k-x}}{\binom{d}{k}}, \quad \text{for } x = 1, 2, \dots, k.$$

1832 *Suppose $k = o(\sqrt{d})$, then for any constant $s > 0$, it holds that $\mathbb{E}[X^s] \simeq k^2/d$. In addition, the*
 1833 *following tail bound holds:*

$$\Pr(X \geq \log k) \lesssim (k^2/d)^{\log k}.$$

1834 *Proof of Lemma J.7.* We first notice that for $x \geq k^2/d$, since $k = o(\sqrt{d})$,

$$\frac{\Pr(X = x+1)}{\Pr(X = x)} = \frac{(k-x)^2}{(x+1)(d-2k-x+1)} = \left(1 - \frac{k}{d}\right)^2 / \left(\frac{d}{k^2} + o\left(\frac{d}{k^2}\right)\right) \lesssim \frac{k^2}{d}. \quad (\text{J.1})$$

1835 This immediately implies the tail bound:

$$\Pr(X \geq \log k) \leq \sum_{j=\lceil \log k \rceil}^k \Pr(X = j) \lesssim \left(\frac{k^2}{d}\right)^{\log k}.$$

1836 Next, for the moment $\mathbb{E}[X^s]$, we first study the magnitude of $\mathbb{P}(X = 0)$, which, by Stirling's
 1837 approximation, is given by

$$\begin{aligned} \Pr(X = 0) &= \frac{((d-k)!)^2}{d! \cdot (d-2k)!} \simeq \frac{(d-k)^{2(d-k)+1} \cdot e^{-2(d-k)}}{d^{d+1/2} \cdot (d-2k)^{(d-2k)+1/2} \cdot e^{-2(d-k)}} \\ &= \frac{(1-2k/d + k^2/d^2)^{d-k+1/2}}{(1-2k/d)^{d-2k+1/2}} \\ &= \left(1 + \frac{1}{(1-2k/d) \cdot d^2/k^2}\right)^{d-k+1/2} \left(1 - \frac{2k}{d}\right)^k. \end{aligned}$$

1838 Since $k = o(\sqrt{d})$, we have $(1 - 2k/d)^k \simeq 1$. Further applying the fact that $(1 + 1/m)^m = \Theta(1)$,

$$\Pr(X = 0) \simeq \left(1 + \frac{1}{(1 - 2k/d) \cdot d^2/k^2}\right)^{(1-2k/d) \cdot d^2/k^2 \cdot k^2/d} = \Theta(1).$$

1839 As a consequence, using the first equality in (J.1), we can lower bound the expectation of X^s by

$$\mathbb{E}[X^s] \geq \Pr(X = 1) = \Pr(X = 0) \cdot \frac{k^2}{d - 2k + 1} \gtrsim \frac{k^2}{d}.$$

1840 For the upper bound, we again use the first equality in (J.1) to get

$$\Pr(X = x + 1) \leq \frac{k^2}{(x + 1)(d - 2k - x + 1)} \cdot \Pr(X = x) \lesssim \frac{\Pr(X = x)}{x + 1} \cdot \frac{k^2}{d}.$$

1841 Recursive application of this inequality yields that $\Pr(X = x) \lesssim \Pr(X = 0) \cdot (k^2/d)^x / x!$, and thus

$$\mathbb{E}[X^s] = \sum_{x=1}^k x^s \cdot \Pr(X = x) \leq \Pr(X = 0) \cdot \sum_{x=1}^k \frac{x^s}{x!} \left(\frac{k^2}{d}\right)^x \lesssim \frac{k^2}{d}.$$

1842 Therefore, we conclude that $\mathbb{E}[X^s] \simeq k^2/d$. This completes the proof. \square