

---

# The Noise Geometry of Stochastic Gradient Descent: A Quantitative and Analytical Characterization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Empirical studies have demonstrated that the noise in stochastic gradient descent  
2 (SGD) aligns favorably with the local geometry of loss landscape. However,  
3 theoretical and quantitative explanations for this phenomenon remain sparse. In this  
4 paper, we offer a comprehensive theoretical investigation into the aforementioned  
5 *noise geometry* for over-parameterized linear (OLMs) models and two-layer neural  
6 networks. We scrutinize both average and directional alignments, paying special  
7 attention to how factors like sample size and input data degeneracy affect the  
8 alignment strength. As a specific application, we leverage our noise geometry  
9 characterizations to study how SGD escapes from sharp minima, revealing that the  
10 escape direction has significant components along flat directions. This is in stark  
11 contrast to GD, which escapes only along the sharpest directions. To substantiate  
12 our theoretical findings, both synthetic and real-world experiments are provided.

## 13 1 Introduction

14 Stochastic gradient descent (SGD) and its variants have become the de facto optimizers for training  
15 machine learning models (Bottou, 1991). Unlike full-batch gradient descent (GD), SGD uses only  
16 mini-batches of data in each iteration, which injects noise into the optimization process. This noise  
17 can have a pronounced impact on both the convergence behavior (Thomas et al., 2020; Wojtowytsch,  
18 2023; Feng and Tu, 2021; Simsekli et al., 2019) and the generalization capabilities (Zhang et al.,  
19 2017; Keskar et al., 2017; Wu et al., 2017; Zhu et al., 2019; Smith et al., 2020) of the algorithm.

20 Zhu et al. (2019); Wu et al. (2020); Xie et al. (2020) showed that SGD noise is highly anisotropic and  
21 in particular, the noise covariance matrix aligns well with the Hessian matrix. As such, they propose  
22 a Hessian-based approximation of the noise covariance:  $\Sigma(\theta) \approx \sigma^2 H(\theta)$ , where  $\Sigma(\theta)$  and  $H(\theta)$   
23 denote the noise covariance and Hessian matrices at  $\theta$ , respectively and  $\sigma$  serves as a small constant  
24 denoting the noise magnitude. Subsequent works (Feng and Tu, 2021; Mori et al., 2022; Wojtowytsch,  
25 2021; Liu et al., 2021) presented an improved Hessian-based approximation:  $\Sigma(\theta) \approx 2L(\theta)H(\theta)$  for  
26 regression problems with square loss, where  $L(\theta)$  denotes the loss value. This refined approximation  
27 acknowledges the fact that the noise magnitude is proportional to the loss value.

28 However, the alignment between SGD noise and local landscape geometry remains empirical observa-  
29 tions, lacking quantitative characterization and theoretical grounding. Hessian-based approximations  
30 are not accurate, as underscored by Thomas et al. (2020). A recent effort by Wu et al. (2022)  
31 employed a normalized cosine similarity between  $\Sigma(\theta)$ —which is close to the Hessian matrix in low  
32 loss regions—and the empirical Fisher matrix  $G(\theta)$  as a metric to quantify the alignment. This metric  
33 is inspired by analyzing the dynamical stability of SGD (Wu et al., 2018) and can be interpreted  
34 as certain type of average alignment. Nevertheless, the analysis in Wu et al. (2022) is restricted to  
35 over-parameterized linear models (OLMs) and operates under the assumption of infinite data, leaving  
36 open questions about the generalizability of such alignment in more practically relevant settings.

37 **Our contribution.** Let  $n, d$  denote the sample size, input dimension, respectively. Then, our  
38 contributions can be summarized as follows.

- 39 • We first extend the average alignment analysis (Wu et al., 2022) to finite sample scenarios,  
40 offering a comprehensive investigation of how factors like sample size and input data degeneracy  
41 impact the alignment strength. We establish that, as long as  $d_{\text{eff}} \gtrsim \log n$ , the alignment strength  
42 is lower-bounded for both OLMs and two-layer neural networks—models not considered in Wu  
43 et al. (2022). Here,  $d_{\text{eff}}$  represents an effective input dimension, and this condition accommodates  
44 the important regimes like  $n \sim \log(d_{\text{eff}})$  (for sparse recovery) and  $n \sim d_{\text{eff}}$  (the proportional  
45 scaling).
- 46 • We then delve into a directional alignment analysis, probing whether the component of noise  
47 energy along a specific direction is proportional to the curvature in that direction. Our results  
48 show that for OLMs, as long as  $n \gtrsim d$ , the strength of directional alignment is lower-bounded  
49 across all directions and the entire parameter space.
- 50 • Lastly, we provide a detailed analysis of the mechanisms by which SGD escapes from sharp  
51 minima by leveraging our noise geometry results. We show that *the escape direction of SGD*  
52 *exhibits significant components along flat directions of the local landscape*. This stands in stark  
53 contrast to GD, which escapes from minima only along the sharpest direction. We also discuss  
54 the implications of this unique escape behavior, providing a preliminary explanation of how  
55 cyclical learning rate (Smith, 2017; Loshchilov and Hutter, 2017) can help find flatter minima.

56 It is worth noting that our theoretical guarantees apply effectively to both isotropic and anisotropic  
57 inputs, and *the guaranteed alignment strength is independent of the degree of overparameterization*.  
58 In addition, all theoretical findings are supported by numerical experiments conducted on both  
59 small-scale and larger-scale models, provided in Appendix C and D. Overall, our work advances  
60 the theoretical understanding of the geometry of SGD noise and provides insights into how SGD  
61 navigates the loss landscape.

## 62 2 Preliminaries

63 Let  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$  be the training set and  $f(\cdot; \boldsymbol{\theta}) : \mathbb{R}^d \rightarrow \mathbb{R}$  be the model parameterized  
64 by  $\boldsymbol{\theta} \in \mathbb{R}^p$ . Let  $\ell_i(\boldsymbol{\theta}) = \frac{1}{2} (f(\mathbf{x}_i; \boldsymbol{\theta}) - y_i)^2$  be the square loss at the  $i$ -th sample and  $\mathcal{L}(\boldsymbol{\theta}) =$   
65  $\frac{1}{n} \sum_{i=1}^n \ell_i(\boldsymbol{\theta})$  be the empirical risk. To minimize  $\mathcal{L}(\cdot)$ , the mini-batch SGD updates as follows  
66  $\boldsymbol{\theta}(t+1) = \boldsymbol{\theta}(t) - \frac{\eta}{B} \sum_{i \in \mathcal{B}_t} \nabla \ell_i(\boldsymbol{\theta}(t))$ , where  $\mathcal{B}_t = \{\gamma_{t,1}, \dots, \gamma_{t,B}\}$  is a batch with size  $|\mathcal{B}_t| = B$ ,  
67 and  $\gamma_{t,1}, \dots, \gamma_{t,B} \stackrel{\text{i.i.d.}}{\sim} \mathbb{U}([n])$ . To isolate the impact of noise, the SGD update is often reformulated  
68 as  $\boldsymbol{\theta}(t+1) = \boldsymbol{\theta}(t) - \eta (\nabla \mathcal{L}(\boldsymbol{\theta}(t)) + \boldsymbol{\xi}(t))$ , where  $\nabla \mathcal{L}(\boldsymbol{\theta}(t))$  is the full-batch gradient and  $\boldsymbol{\xi}(t)$   
69 represents the mini-batch noise satisfying  $\mathbb{E}[\boldsymbol{\xi}(t)] = 0$ ,  $\mathbb{E}[\boldsymbol{\xi}(t)\boldsymbol{\xi}(t)^\top] = \Sigma(\boldsymbol{\theta}(t))/B$  with the noise  
70 covariance given by  $\Sigma(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(\boldsymbol{\theta}) \nabla \ell_i(\boldsymbol{\theta})^\top - \nabla \mathcal{L}(\boldsymbol{\theta}) \nabla \mathcal{L}(\boldsymbol{\theta})^\top$ . In the above setup, the  
71 Hessian matrix of the empirical risk is given by  $H(\boldsymbol{\theta}) = G(\boldsymbol{\theta}) + \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i; \boldsymbol{\theta}) - y_i) \nabla^2 f(\mathbf{x}_i; \boldsymbol{\theta})$ ,  
72 where  $G(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{x}_i; \boldsymbol{\theta}) \nabla f(\mathbf{x}_i; \boldsymbol{\theta})^\top$  is the empirical Fisher matrix. When the fit errors are  
73 small, we have  $G(\boldsymbol{\theta}) \approx H(\boldsymbol{\theta})$  and in particular, for global minima  $\boldsymbol{\theta}^*$ ,  $H(\boldsymbol{\theta}^*) = G(\boldsymbol{\theta}^*)$ . Additionally,  
74 for linear regression  $f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}$ ,  $H(\boldsymbol{\theta}) = G(\boldsymbol{\theta}) \equiv \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ .

75 **Over-parameterized linear models (OLMs).** An OLM is defined as  $f(\mathbf{x}; \boldsymbol{\theta}) = F(\boldsymbol{\theta})^\top \mathbf{x}$ , where  
76  $F : \mathbb{R}^p \rightarrow \mathbb{R}^d$  denotes a general re-parameterization function. Although  $f(\cdot; \boldsymbol{\theta})$  only represents linear  
77 functions, the corresponding loss landscape can be highly non-convex. Some typical examples include  
78 (i) the linear model  $F(\mathbf{w}) = \mathbf{w}$ ; (ii) the diagonal linear network:  $F(\boldsymbol{\theta}) = (\alpha_1^2 - \beta_1^2, \dots, \alpha_d^2 - \beta_d^2)^\top$ ;  
79 and (iii) the linear network:  $F(\boldsymbol{\theta}) = W_1 W_2 \dots W_L$ . Notably, OLMs have been widely used to  
80 analyze the optimization and implicit bias of SGD (Arora et al., 2019; Woodworth et al., 2020; Pesme  
81 et al., 2021; HaoChen et al., 2021; Azulay et al., 2021).

82 **Noise Geometry.** Before proceeding to our refined characterization of the noise geometry, we first  
83 recall two existing results on quantifying the geometry of SGD noise.

- 84 • Mori et al. (2022) proposed the following Hessian-based approximation:  $\Sigma(\boldsymbol{\theta}) \approx 2\mathcal{L}(\boldsymbol{\theta})G(\boldsymbol{\theta})$ .  
85 It reveals 1) the noise magnitude is proportional to the loss value; 2) the noise covariance aligns  
86 with the Fisher matrix. This approximation is intuitive and helpful for understanding, but it  
87 cannot be accurate in general.
- 88 • *Online SGD for OLMs with Gaussian inputs.* Suppose  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, S)$  and  $n = \infty$  (i.e.,  
89 online SGD). For OLMs, Wu et al. (2022) derived the following analytical expression  
90  $\Sigma(\boldsymbol{\theta}) = 2\mathcal{L}(\boldsymbol{\theta})G(\boldsymbol{\theta}) + \nabla \mathcal{L}(\boldsymbol{\theta}) \nabla \mathcal{L}(\boldsymbol{\theta})^\top$ .

### 91 3 Average Alignment

92 Let  $\Sigma_1(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(\boldsymbol{\theta}) \nabla \ell_i(\boldsymbol{\theta})^\top$ ,  $\Sigma_2(\boldsymbol{\theta}) = \nabla \mathcal{L}(\boldsymbol{\theta}) \nabla \mathcal{L}(\boldsymbol{\theta})^\top$ . Then  $\Sigma(\boldsymbol{\theta}) = \Sigma_1(\boldsymbol{\theta}) - \Sigma_2(\boldsymbol{\theta})$ . It  
 93 is commonly believed that the magnitude of the full-batch gradient  $\nabla \mathcal{L}$  is relatively small compared  
 94 to the sample gradients  $\{\nabla \ell_i\}_i$ . Consequently, the influence of  $\Sigma_2(\boldsymbol{\theta})$  would be negligible compared  
 95 to  $\Sigma_1(\boldsymbol{\theta})$ . Following Wu et al. (2022), we consider the following metrics of quantifying average  
 96 alignment:  $\mu(\boldsymbol{\theta}) = \frac{\text{Tr}(\Sigma(\boldsymbol{\theta})G(\boldsymbol{\theta}))}{2\mathcal{L}(\boldsymbol{\theta})\|G(\boldsymbol{\theta})\|_F^2}$ .

97 Wu et al. (2022) guarantees  $\mu(\boldsymbol{\theta}) \geq 1$  in an infinite data scenario. The following theorem extends  
 98 it to finite-sample cases and the proof can be found in Appendix G. To simplify the statement, we  
 99 define the effective dimension of inputs as  $d_{\text{eff}} := \min\{\text{srk}(S), \text{srk}(S^2)\}$ , where  $S$  represents the  
 100 input covariance matrix and  $\text{srk}(S) = \text{tr}(S)/\|S\|_2$  is the stable rank of  $S$ . In particular, when  $S$  is  
 101 isotropic, we have  $d_{\text{eff}} = d$ .

102 **Theorem 3.1** (OLM). *Consider OLMs and assume  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, S)$ . For any  $\epsilon, \delta \in (0, 1)$ ,*  
 103 (a) *if  $n/\log(n/\delta) \gtrsim 1/\epsilon^2$  and  $d_{\text{eff}} \gtrsim \log(n/\delta)/\epsilon^2$ , then w.p. at least  $1 - \delta$ , it holds that*  
 104 
$$\inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \mu(\boldsymbol{\theta}) \geq \frac{(1-\epsilon)^2}{(1+\epsilon)^2 \text{cond}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top)};$$

105 (b) *if  $n \gtrsim d + \log(1/\delta)$ , then w.p. at least  $1 - \delta$ , it holds that  $\inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \mu(\boldsymbol{\theta}) \gtrsim 1$ .*

106 Result (a) is established by leveraging the high dimensionality of inputs, as stated by the condition  
 107  $d_{\text{eff}} \gtrsim \log n$ , which is particularly relevant for low-sample regimes. Notably, this includes the  
 108 important regimes like  $n \sim \log(d_{\text{eff}})$  (for sparse recovery) and  $n \sim d_{\text{eff}}$  (the proportional scaling).  
 109 In contrast, result (b) is pertinent to the enough-data regime where  $n \gtrsim d$ . Notably, the alignment  
 110 holds no matter how degenerate the covariance matrix is. In a summary, these two results are  
 111 complementary and collectively span all the regimes of interest.

112 **Example.** Consider the isotropic case where  $S = I_d$  and linear regression  $F(\mathbf{w}) = \mathbf{w}$ . In this case,  
 113  $\nabla F(\mathbf{w}) \equiv I_d$  and thus, Theorem 3.1 implies that it holds that  $\inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \mu(\boldsymbol{\theta}) \gtrsim 1$  as long as  $n \gtrsim 1$ .

114 Consider two-layer neural networks given by  $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^m a_k \phi(\mathbf{b}_k^\top \mathbf{x})$  with  $a_k \in \{\pm 1\}$  to be  
 115 fixed. We use  $\boldsymbol{\theta} = (\mathbf{b}_1^\top, \dots, \mathbf{b}_m^\top)^\top \in \mathbb{R}^{md}$  to denote the concatenation of all trainable parameters.  
 116 Here,  $\phi: \mathbb{R} \mapsto \mathbb{R}$  is an activation function with a nondegenerate derivative as defined below.

117 **Assumption 3.2.** There exist constants  $\beta > \alpha > 0$  such that  $\alpha \leq \phi'(z) \leq \beta$  holds for any  $z \in \mathbb{R}$ .

118 **Example 3.3.** (i) A typical activation function that satisfies Assumption 3.2 is  $\alpha$ -Leaky ReLU:  $\phi(z) =$   
 119  $\max\{\alpha z, z\}$ , where  $\alpha \in (0, 1)$ . (ii) Moreover, the assumption also holds for Sigmoid with the trunca-  
 120 tion trick (to prevent gradient vanishing of Sigmoid):  $\phi(z) = 1/(1 + \exp(-\text{sgn}(z) \min\{|z|, M\}))$ ,  
 121 where  $M > 0$  is the truncation constant.

122 **Theorem 3.4** (2NN). *Consider the two-layer network  $f(\cdot; \boldsymbol{\theta})$  with the activation function satisfying*  
 123 *Assumption 3.2 and assume  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, S)$ . For any  $\epsilon, \delta \in (0, 1)$ , if  $n/\log(n/\delta) \gtrsim 1/\epsilon^2$*   
 124 *and  $d_{\text{eff}} \gtrsim \log(n/\delta)/\epsilon^2$ , then w.p. at least  $1 - \delta$ , it holds that  $\inf_{\boldsymbol{\theta} \in \mathbb{R}^{md}} \mu(\boldsymbol{\theta}) \geq \frac{\alpha^2(1-\epsilon)^2}{\beta^2(1+\epsilon)^2}$ .*

125 **Remark 3.5.** We would like to emphasize that the conditions presented in Theorem 3.1 and 3.4 are  
 126 independent of the model size  $p$ .

127 The numerical validation is referred to Appendix C.

### 128 4 Directional Alignment

129 In Section 3, we focused solely on average alignment. Subsequently, we delve into a specific type of  
 130 directional alignment: *whether noise energy along a direction is proportional to the curvature of loss*  
 131 *landscape along that direction.* To this end, we define the following metric to measure the strength of  
 132 directional alignment.

133 **Definition 4.1** (Directional Alignment). Given  $\mathbf{v} \in \mathbb{R}^p$ , the alignment along  $\mathbf{v}$  is defined as  $g(\boldsymbol{\theta}; \mathbf{v}) :=$   
 134  $\frac{\mathbf{v}^\top \Sigma(\boldsymbol{\theta}) \mathbf{v}}{2\mathcal{L}(\boldsymbol{\theta})(\mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v})}$ , where  $\mathbf{v}^\top \Sigma(\boldsymbol{\theta}) \mathbf{v} = \mathbb{E}[(\boldsymbol{\xi}(\boldsymbol{\theta})^\top \mathbf{v})^2]$  denotes the noise energy along direction  $\mathbf{v}$ ,  
 135  $\mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}$  is the curvature of loss landscape along  $\mathbf{v}$ , and  $2\mathcal{L}(\boldsymbol{\theta})$  is only a scaling factor.

136 **Theorem 4.2** (One-sided bound). *Consider OLMs and assume  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, S)$ . For*  
 137 *any  $\delta \in (0, 1)$ , if  $n \gtrsim d + \log(1/\delta)$ , then w.p. at least  $1 - \delta$ , we have  $\inf_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p} g(\boldsymbol{\theta}; \mathbf{v}) \gtrsim 1$ .*

138 This theorem establishes that a sample size satisfying  $n \gtrsim d$  is sufficient to guarantee a uniform lower  
 139 bound for alignment across all directions and the entire parameter space. The subsequent theorem  
 140 builds upon this by offering a two-sided bound on alignment strength, albeit at the cost of requiring a  
 141 larger sample size.

142 **Theorem 4.3** (Two-sided bound). *Consider OLMs and assume  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, S)$ .*  
 143 *For any  $\epsilon, \delta \in (0, 1)$ , if  $n \gtrsim \max\{(d^2 \log^2(1/\epsilon) + \log^2(1/\delta)) / \epsilon, (d \log(1/\epsilon) + \log(1/\delta)) / \epsilon^2\}$ ,*  
 144 *then w.p. at least  $1 - \delta$ , we have the following two-side uniform bounds for the directional*  
 145 *alignment: (i).  $\frac{1-\epsilon}{(1+\epsilon)^2} \leq \inf_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p} g(\boldsymbol{\theta}; \mathbf{v}) \leq \sup_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p} g(\boldsymbol{\theta}; \mathbf{v}) \leq \frac{2+\epsilon}{(1-\epsilon)^2}$ ; (ii).  $\frac{1-\epsilon}{(1+\epsilon)^2} \leq$*   
 146  *$\inf_{\boldsymbol{\theta} \in \mathbb{R}^p, \langle \mathbf{v}, \nabla \mathcal{L}(\boldsymbol{\theta}) \rangle = 0} g(\boldsymbol{\theta}; \mathbf{v}) \leq \sup_{\boldsymbol{\theta} \in \mathbb{R}^p, \langle \mathbf{v}, \nabla \mathcal{L}(\boldsymbol{\theta}) \rangle = 0} g(\boldsymbol{\theta}; \mathbf{v}) \leq \frac{1+\epsilon}{(1-\epsilon)^2}$ .*

147 Notably, for directions satisfying  $\mathbf{v} \perp \nabla \mathcal{L}(\boldsymbol{\theta})$ , the alignment strength is nearly 1. The proofs of the  
 148 above two theorems are deferred to Appendix H. The numerical validation is referred to Appendix C.

## 149 5 How SGD Escapes from Sharp Minima

150 Let  $\mathbf{w} = \boldsymbol{\theta} - \boldsymbol{\theta}^*$  and  $\mathbf{z}_i = \nabla f(\mathbf{x}_i; \boldsymbol{\theta}^*)$ . Then,  $G(\boldsymbol{\theta}^*) = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^\top$  and the linearized SGD  
 151 of iterates as follows  $\mathbf{w}(t+1) = \mathbf{w}(t) - \eta(G(\boldsymbol{\theta}^*)\mathbf{w}(t) + \boldsymbol{\xi}(t))$ , where  $\boldsymbol{\xi}(t)$  is the SGD noise. In  
 152 addition, in this section, we simply use  $\mathcal{L}(\mathbf{w}) = \frac{1}{2} \mathbf{w}^\top G(\boldsymbol{\theta}^*) \mathbf{w}$  to denote the corresponding loss. We  
 153 make the following assumption on the noise alignment.

154 **Assumption 5.1** (Eigen-directional alignment). let  $G(\boldsymbol{\theta}^*) = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$  be the eigen decom-  
 155 position of  $G(\boldsymbol{\theta}^*)$ . Assume that there exist  $A_1, A_2 > 0$  such that it holds for any  $\mathbf{w} \in \mathbb{R}^d$ ,  
 156  $A_1 \mathcal{L}(\mathbf{w}) \lambda_i \leq \mathbb{E}[|\boldsymbol{\xi}(\mathbf{w})^\top \mathbf{u}_i|^2] \leq A_2 \mathcal{L}(\mathbf{w}) \lambda_i$ .

157 For linear models under the setting of Theorem 4.3, Assumption 5.1 is provably valid. It is important  
 158 to clarify, however, that the above assumption only requires the alignment along eigen-directions,  
 159 which is considerably less stringent compared to the uniform directional alignment specified in  
 160 Theorem 4.3.

161 **Eigen-decomposition of SGD.** By leveraging Assumption 5.1, we can analyze the SGD dynamics  
 162 in the eigenspace. Let  $\mathbf{w}(t) = \sum_{i=1}^d w_i(t) \mathbf{u}_i$  with  $w_i(t) = \mathbf{u}_i^\top \mathbf{w}(t)$ . Then,  $w_i(t+1) = (1 -$   
 163  $\eta \lambda_i) w_i(t) + \eta \boldsymbol{\xi}(t)^\top \mathbf{u}_i$ . Taking the expectation of the square of both sides, we obtain  $\mathbb{E}[w_i^2(t+1)] =$   
 164  $(1 - \eta \lambda_i)^2 \mathbb{E}[w_i^2(t)] + \eta^2 \mathbb{E}[|\mathbf{u}_i^\top \boldsymbol{\xi}(t)|^2]$ , where the noise term:  $\mathbb{E}[|\mathbf{u}_i^\top \boldsymbol{\xi}(t)|^2] \sim \lambda_i \mathcal{L}(\mathbf{w}_t)$  according  
 165 to Assumption 5.1.

166 Let  $X_t = \sum_{i=1}^k \lambda_i \mathbb{E}[w_i^2(t)]$ ,  $Y_t = \sum_{i=k+1}^d \lambda_i \mathbb{E}[w_i^2(t)]$ , denoting the components of loss energy  
 167 along sharp and flat directions, respectively. Let  $D_k(t) = Y_t / X_t$ , which measures the concentration  
 168 of loss energy along flat directions. Analogously, let  $P_k(t) = \sum_{i=k+1}^d \mathbb{E}[w_i^2(t)] / \sum_{i=1}^k \mathbb{E}[w_i^2(t)]$ ,  
 169 which measure the concentration of variance along flat directions. It is easy to show that  $P_k(t) \geq$   
 170  $D_k(t) \lambda_k / \lambda_{k+1}$ . Therefore, when  $\lambda_k / \lambda_{k+1}$  is lower bounded, a concentration of loss energy along  
 171 flat directions can lead to a similar concentration in terms of variance.

172 **Theorem 5.2** (Escape of SGD). *Suppose Assumption 5.1 holds and let  $\eta = \frac{\beta}{\|G(\boldsymbol{\theta}^*)\|_F}$ . Then, there*  
 173 *exists absolute constants  $c_1, c_2 > 0$  such that if  $\beta \geq c_1$ , then SGD will escape from that minima and*  
 174 *for any  $k \in [d]$ , it holds that when  $t \geq \max\left\{1, \frac{\log(c_2 / \eta (\sum_{i=1}^k \lambda_i^2)^{1/2})}{\log \beta}\right\}$ :  $D_k(t) \gtrsim \frac{\sum_{i=k+1}^d \lambda_i^2}{\sum_{i=1}^k \lambda_i^2}$ .*

175 The proof can be found in Appendix I. This theorem reveals that during SGD's escape process, the  
 176 loss rapidly accumulates a significant component along flat directions of the loss landscape. The  
 177 precise loss ratio between the flat and sharp directions is governed by the spectrum of Hessian matrix.  
 178 In particular,  $D_1(t) \gtrsim \text{srk}(G^2) - 1$ , indicating that in high dimension, i.e.,  $\text{srk}(G^2) \gg 1$ , the loss  
 179 energy along the sharpest directions becomes negligible during the SGD's escape process. This  
 180 stands in stark contrast to GD, which always escapes along the sharpest direction:

181 **Proposition 5.3** (Escape of GD). *Consider GD with learning rate  $\eta = \beta / \lambda_1$ . If  $\beta > 2$ , then*  
 182  $D_1(t) \leq \sum_{i=2}^d \frac{\lambda_i (1 - \eta \lambda_i)^{2t} w_i^2(0)}{\lambda_1 (1 - \eta \lambda_1)^{2t} w_1^2(0)}$ .

183 In particular, if  $w_1(0) \neq 0$  and  $\lambda_1 > \lambda_2$ , then the above proposition implies that  $D_1(t)$  decreases to 0  
 184 exponentially fast for GD. The numerical validation is referred to Appendix C.

185 Furthermore, as an implication of SGD's escaping direction, we explain the **implicit bias of cyclical**  
 186 **learning rate** in Appendix B.

## 187 References

- 188 Kwangjun Ahn, Jingzhao Zhang, and Suvrit Sra. Understanding the unstable convergence of gradient  
189 descent. In *International Conference on Machine Learning*, pages 247–257. PMLR, 2022.
- 190 Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix  
191 factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- 192 Shahar Azulay, Edward Moroshko, Mor Shpigel Nacson, Blake E Woodworth, Nathan Srebro, Amir  
193 Globerson, and Daniel Soudry. On the implicit bias of initialization shape: Beyond infinitesimal  
194 mirror descent. In *International Conference on Machine Learning*, pages 468–477. PMLR, 2021.
- 195 Léon Bottou. Stochastic gradient learning in neural networks. *Proceedings of Neuro-Nimes*, 91(8),  
196 1991.
- 197 Jian-Feng Cai, Meng Huang, Dong Li, and Yang Wang. Nearly optimal bounds for the global  
198 geometric landscape of phase retrieval. *arXiv preprint arXiv:2204.09416*, 2022.
- 199 Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ammeet Talwalkar. Gradient descent on  
200 neural networks typically occurs at the edge of stability. In *International Conference on Learning  
201 Representations*, 2020.
- 202 Hadi Daneshmand, Jonas Kohler, Aurelien Lucchi, and Thomas Hofmann. Escaping saddles with  
203 stochastic gradients. In *International Conference on Machine Learning*, pages 1155–1164. PMLR,  
204 2018.
- 205 Yu Feng and Yuhai Tu. The inverse variance–flatness relation in stochastic gradient descent is critical  
206 for finding flat minima. *Proceedings of the National Academy of Sciences*, 118(9), 2021.
- 207 Jonas Geiping, Micah Goldblum, Phillip E Pope, Michael Moeller, and Tom Goldstein. Stochastic  
208 training is not necessary for generalization. *arXiv preprint arXiv:2109.14119*, 2021.
- 209 Botao Hao, Yasin Abbasi Yadkori, Zheng Wen, and Guang Cheng. Bootstrapping upper confidence  
210 bound. *Advances in neural information processing systems*, 32, 2019.
- 211 Jeff Z HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. Shape matters: Understanding the implicit  
212 bias of the noise covariance. In *Conference on Learning Theory*, pages 2315–2357. PMLR, 2021.
- 213 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
214 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
215 pages 770–778, 2016.
- 216 S. Hochreiter and J. Schmidhuber. Flat minima. *Neural Computation*, 9(1):1–42, 1997.
- 217 Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E Hopcroft, and Kilian Q Weinberger. Snap-  
218 shot ensembles: Train 1, get  $M$  for free. In *International Conference on Learning Representations*,  
219 2018.
- 220 N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. On large-batch training for  
221 deep learning: Generalization gap and sharp minima. In *International Conference on Learning  
222 Representations (ICLR)*, 2017.
- 223 Bobby Kleinberg, Yuanzhi Li, and Yang Yuan. An alternative view: When does SGD escape local  
224 minima? In *International Conference on Machine Learning*, pages 2698–2707. PMLR, 2018.
- 225 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images, 2009.  
226 URL <https://www.cs.toronto.edu/~kriz/cifar.html>.
- 227 Sungyoon Lee and Cheongjae Jang. A new characterization of the edge of stability based on a  
228 sharpness measure aware of batch gradient distribution. In *The Eleventh International Conference  
229 on Learning Representations*, 2022.
- 230 Kangqiao Liu, Liu Ziyin, and Masahito Ueda. Noise and fluctuation of finite learning rate stochastic  
231 gradient descent. In *International Conference on Machine Learning*, pages 7045–7056. PMLR,  
232 2021.



- 233 Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In  
234 *International Conference on Learning Representations*, 2017.
- 235 Chao Ma, Daniel Kunin, Lei Wu, and Lexing Ying. Beyond the quadratic approximation: The  
236 multiscale structure of neural network loss landscapes. *Journal of Machine Learning*, 1(3):  
237 247–267, 2022.
- 238 Takashi Mori, Liu Ziyin, Kangqiao Liu, and Masahito Ueda. Power-law escape rate of SGD. In  
239 *International Conference on Machine Learning*, pages 15959–15975. PMLR, 2022.
- 240 Scott Pesme, Loucas Pillaud-Vivien, and Nicolas Flammarion. Implicit bias of SGD for diagonal  
241 linear networks: a provable benefit of stochasticity. *Advances in Neural Information Processing*  
242 *Systems*, 34, 2021.
- 243 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image  
244 recognition. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- 245 Umut Simsekli, Levent Sagun, and Mert Gurbuzbalaban. A tail-index analysis of stochastic gradient  
246 noise in deep neural networks. In *International Conference on Machine Learning*, pages 5827–  
247 5837. PMLR, 2019.
- 248 Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference*  
249 *on applications of computer vision (WACV)*, pages 464–472. IEEE, 2017.
- 250 Samuel Smith, Erich Elsen, and Soham De. On the generalization benefit of noise in stochastic  
251 gradient descent. In *International Conference on Machine Learning*, pages 9058–9067. PMLR,  
252 2020.
- 253 Valentin Thomas, Fabian Pedregosa, Bart Merriënboer, Pierre-Antoine Manzagol, Yoshua Bengio,  
254 and Nicolas Le Roux. On the interplay between noise and curvature and its effect on optimization  
255 and generalization. In *International Conference on Artificial Intelligence and Statistics*, pages  
256 3503–3513. PMLR, 2020.
- 257 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,  
258 volume 47. Cambridge university press, 2018.
- 259 Stephan Wojtowytsch. Stochastic gradient descent with noise of machine learning type. part II:  
260 Continuous time analysis. *arXiv preprint arXiv:2106.02588*, 2021.
- 261 Stephan Wojtowytsch. Stochastic gradient descent with noise of machine learning type part i: Discrete  
262 time analysis. *Journal of Nonlinear Science*, 33(3):45, 2023.
- 263 Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan,  
264 Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparameterized models. In  
265 *Conference on Learning Theory*, pages 3635–3673. PMLR, 2020.
- 266 Jingfeng Wu, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu. On the  
267 noisy gradient descent that generalizes as SGD. In *International Conference on Machine Learning*,  
268 pages 10367–10376. PMLR, 2020.
- 269 Lei Wu and Weijie J Su. The implicit regularization of dynamical stability in stochastic gradient  
270 descent. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202,  
271 pages 37656–37684. PMLR, 23–29 Jul 2023.
- 272 Lei Wu, Zhanxing Zhu, and Weinan E. Towards understanding generalization of deep learning:  
273 Perspective of loss landscapes. *arXiv preprint arXiv:1706.10239*, 2017.
- 274 Lei Wu, Chao Ma, and Weinan E. How SGD selects the global minima in over-parameterized  
275 learning: A dynamical stability perspective. *Advances in Neural Information Processing Systems*,  
276 31:8279–8288, 2018.
- 277 Lei Wu, Mingze Wang, and Weijie J Su. The alignment property of SGD noise and how it helps  
278 select flat minima: A stability analysis. *Advances in Neural Information Processing Systems*, 35:  
279 4680–4693, 2022.

- 280 Zeke Xie, Issei Sato, and Masashi Sugiyama. A diffusion theory for deep learning dynamics:  
281 Stochastic gradient descent exponentially favors flat minima. In *International Conference on*  
282 *Learning Representations*, 2020.
- 283 Zeke Xie, Xinrui Wang, Huishuai Zhang, Issei Sato, and Masashi Sugiyama. Adaptive inertia:  
284 Disentangling the effects of adaptive learning rate and momentum. In *International conference on*  
285 *machine learning*, pages 24430–24459. PMLR, 2022.
- 286 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understand-  
287 ing deep learning requires rethinking generalization. In *International Conference on Learning*  
288 *Representations*, 2017.
- 289 Hongyi Zhang, Yann N. Dauphin, and Tengyu Ma. Residual learning without normalization via better  
290 initialization. In *International Conference on Learning Representations*, 2019.
- 291 Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, and Weinan E. Towards  
292 theoretically understanding why SGD generalizes better than Adam in deep learning. *Advances in*  
293 *Neural Information Processing Systems*, 33, 2020.
- 294 Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. The anisotropic noise in stochastic  
295 gradient descent: Its behavior of escaping from sharp minima and regularization effects. In  
296 *International Conference on Machine Learning*, pages 7654–7663. PMLR, 2019.
- 297 Liu Ziyin, Kangqiao Liu, Takashi Mori, and Masahito Ueda. Strength of minibatch noise in SGD. In  
298 *International Conference on Learning Representations*, 2022.

299  
300  
301

---

# Appendix

---

302	<b>A Other Related Work</b>	<b>8</b>
303	<b>B Explaining the Implicit Bias of Cyclical Learning Rate</b>	<b>9</b>
304	<b>C Small-scale Experiments</b>	<b>9</b>
305	C.1 Average Alignment . . . . .	9
306	C.2 Directional Alignment . . . . .	10
307	C.3 Escaping Direction . . . . .	10
308	<b>D Larger-scale Experiments for Deep Neural Networks</b>	<b>11</b>
309	<b>E Conclusion and Future Work</b>	<b>12</b>
310	<b>F Experimental Setups</b>	<b>12</b>
311	<b>G Proofs in Section 3: Average alignment</b>	<b>13</b>
312	G.1 Proof of Theorem 3.1 (a) . . . . .	13
313	G.2 Proof of Theorem 3.1 (b) . . . . .	17
314	G.3 Proof of Theorem 3.4 . . . . .	17
315	<b>H Proofs in Section 4: Directional Alignment</b>	<b>18</b>
316	H.1 Proof of Theorem 4.2 . . . . .	19
317	H.2 Proof of Theorem 4.3 . . . . .	20
318	<b>I Proofs in Section 5: Escape directions</b>	<b>27</b>
319	I.1 Proof of Theorem 5.2 . . . . .	27
320	I.2 Proof of Proposition 5.3 . . . . .	29
321	<b>J Useful Inequalities</b>	<b>29</b>

## 322 A Other Related Work

323 **Noise geometry.** [Ziyin et al. \(2022\)](#) provides a detailed analysis of the noise structure of online  
324 SGD for linear regression. We instead consider nonlinear models and finite-sample regimes. We also  
325 acknowledge the existence of works such as [Simsekli et al. \(2019\)](#); [Zhou et al. \(2020\)](#), which argue  
326 that the magnitude of SGD noise is heavy-tailed. However, our particular focus is on the noise shape  
327 and the observation that the noise magnitude is directly proportional to the loss value.

328 **Escape from minima and saddle points** The phenomenon of SGD escaping from sharp minima  
329 exponentially fast was initially studied in [Zhu et al. \(2019\)](#) as an indicator of how much SGD dislikes  
330 sharp minima. This provides an explanation of the famous “flat minima hypothesis” ([Hochreiter and  
331 Schmidhuber, 1997](#); [Keskar et al., 2017](#); [Wu and Su, 2023](#))—one of the most important observations in  
332 explaining the implicit regularization of SGD. However, existing analyses of the escape phenomenon  
333 have primarily focused on the escape rate ([Wu et al., 2018](#); [Zhu et al., 2019](#); [Xie et al., 2020](#); [Mori  
334 et al., 2022](#); [Ziyin et al., 2022](#)). In contrast, we extend this focus by providing analysis of escape



335 direction, which is enabled by our characterizations of the noise geometry. Kleinberg et al. (2018)  
 336 introduced an alternative perspective, positing that SGD circumvents local minima by navigating  
 337 an effective loss landscape that results from the convolution of the original landscape with SGD  
 338 noise. In this context, our noise geometry characterizations can be beneficial in understanding the  
 339 effective loss landscape. In addition, prior works like (Daneshmand et al., 2018; Xie et al., 2022) has  
 340 illustrated that the alignment of noise with local geometry facilitates the rapid saddle-point escape of  
 341 SGD. Our work offers theoretical substantiation for the alignment assumptions in these studies.

## 342 B Explaining the Implicit Bias of Cyclical Learning Rate

343 Gaining insights into the escape direction of SGD can be valuable for understanding its optimization  
 344 dynamics, generalization properties, and the overall behavior. A more detailed discussion on this  
 345 topic is available in Section E. In this section, however, we concentrate a specific example, illustrating  
 346 the role of escape direction in enhancing the implicit bias of SGD through Cyclical Learning Rate  
 347 (CLR) (Smith, 2017; Loshchilov and Hutter, 2017). As shown in Figure 2 of Huang et al. (2018),  
 348 utilizing CLR enables SGD to cyclically escapes from (when increasing LR) and slides into (when  
 349 decreasing LR) sharp regions, ultimately progressing towards flatter minima. We hypothesize that  
 350 escape along flat directions plays a pivotal role in guiding SGD towards flatter region in this process.

351 Following Ma et al. (2022), we consider a toy OLM  
 352  $f(x; \mathbf{w}) = (w_2/\sqrt{w_1^2 + 1})x$  with  $x \sim \mathcal{N}(0, 1)$ . For sim-  
 353 plicity, we consider the online setting, where the landscape

$$\mathcal{L}(\mathbf{w}) = w_2^2/[2(w_1^2 + 1)].$$

355 The global minima valley is  $S = \{\mathbf{w} : w_2 = 0\}$  and for  
 356  $\mathbf{w} \in S$ ,  $\text{tr}[\nabla^2 \mathcal{L}(\mathbf{w})] = 1/(1 + w_1^2)$ . Hence, the minimum  
 357 gets flatter along the valley  $S$  when  $|w_1|$  grows up. In  
 358 Figure 1, we visualize the trajectories for both SGD+CLR  
 359 and GD+CLR. One can observe that

- 360 • SGD escape from the minima along both the flat  
 361 direction  $e_1$  and sharp direction  $e_2$ . The component  
 362 of along  $e_1$  leads to considerable increase in  $w_1^2(t)$ ,  
 363 facilitating the movement towards flatter region along  
 364 the minimum valley  $S$ .
- 365 • On the contrary, GD escapes only along  $e_2$ , yielding  
 366 no increase in  $w_1^2(t)$ . Thus, we cannot observe clear  
 367 movement towards flatter region for GD+CLR.

368 Thus, in this toy model, the fact that SGD escapes along flat directions is crucial in amplifying the  
 369 implicit bias towards flat minima.

370 Nonetheless, understanding how the above mechanism manifests in practice remains an open question  
 371 that warrants further investigation. We defer this topic to future work, as the primary focus of this  
 372 paper is to understand the noise geometry rather than exhaustively explore its applications.

## 373 C Small-scale Experiments

### 374 C.1 Average Alignment

375 In this section, we present small-scale experiments to corroborate our theoretical results with a 4-layer  
 376 linear network and two-layer ReLU network (both layers are trainable). Both isotropic and anisotropic  
 377 input distributions are examined and in particular, for the anisotropic case, we set  $\lambda_k^2(S) = 1/\sqrt{k}$ . As  
 378 for sample size, we set  $n = 5 \log(d_{\text{eff}})$  to focus on the low-sample regime. The results are reported  
 379 in Figure 2 and it is evident that across all examined scenarios, the alignment strength is consistently  
 380 lower-bounded and independent of the model size.

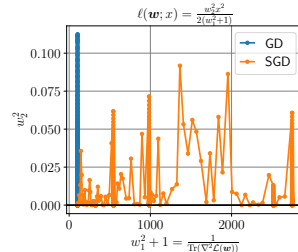


Figure 1: Visualization of the trajectories of SGD+CLR v.s. GD+CLR for our toy model. Both cases use the same CLR schedule. We can observe that SGD+CLR moves significantly towards flatter region, while GD+CLR only oscillates along the sharpest direction. We have extensively tuned the learning rates for GD+CLR but do not observe significant movement towards flatter region in any case.

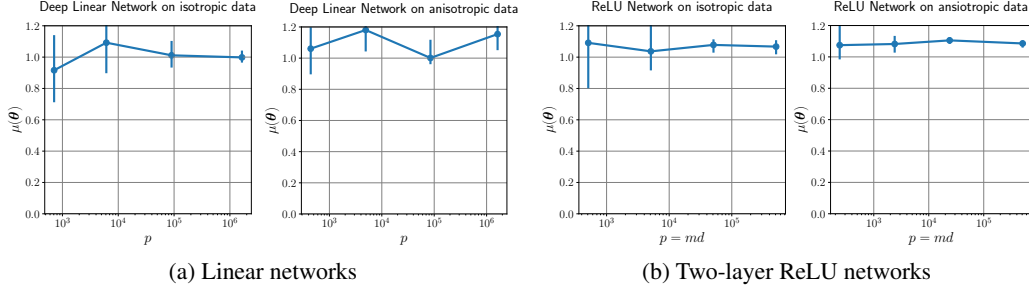


Figure 2: The alignment strength is independent of model size. Two types of models: 4-layer linear network, and two-layer neural network are examined. In experiments, we set  $n = 5 \log(d_{\text{eff}})$ ,  $d_{\text{eff}} = 50$ . The error bar corresponds to the standard deviation over 20 independent runs.

### 381 C.2 Directional Alignment

382 In this experiment, we consider the alignment along the eigen-directions of Hessian matrix. Let  
 383  $G(\theta) = \sum_k \lambda_k(\theta) \mathbf{u}_k(\theta) \mathbf{u}_k(\theta)^\top$  be the eigen-decomposition of  $G(\theta)$  respectively, where  $\{\lambda_k(\theta)\}_k$   
 384 are the eigenvalues in a decreasing order and  $\{\mathbf{u}_k(\theta)\}$  are the corresponding eigen-directions. Note  
 385 that  $\lambda_k(\theta)$  is the curvature of local landscape along  $\mathbf{u}_k(\theta)$ . Decompose SGD noise along these eigen-  
 386 directions:  $\xi(\theta) = \sum_k r_k(\theta) \mathbf{u}_k(\theta)$ , where  $r_k(\theta) = \xi(\theta)^\top \mathbf{u}_k(\theta)$  denotes the noise component in  
 387 the direction of  $\mathbf{u}_k(\theta)$ . Consequently, the (scaled) expected noise magnitude in the direction  $\mathbf{u}_k(\theta)$  is  
 388 given by  $\alpha_k(\theta) = \mathbb{E}[r_k^2(\theta)]/2\mathcal{L}(\theta) = \mathbf{u}_k^\top \Sigma(\theta) \mathbf{u}_k / 2\mathcal{L}(\theta)$ . For comparison, let  $\{\mu_k(\theta)\}_k$  denote  
 389 the eigenvalues of  $\Sigma(\theta)/2\mathcal{L}(\theta)$ . When clear from the context, we will omit dependence on  $\theta$  for  
 390 simplicity.

391 In Figure 3a, we examine linear regression in the regimes with limited data. Surprisingly, even  
 392 with significantly fewer samples, we still observed that the noise energy along each eigen-direction  
 393 remained roughly proportional to the corresponding curvature and the ratio is close 1. However,  
 394 we noticed that the eigenvalues of  $\Sigma(\theta)/2\mathcal{L}(\theta)$  decayed much faster than that of  $G(\theta)$ , indicating  
 395 that the condition  $n \gtrsim d$  stated in Theorem 4.2 is necessary to ensure uniform alignment across all  
 396 directions. In Figure 3b, we further consider the classification of CIFAR-10 with a small convolutional  
 397 neural network (CNN) and fully-connected neural network (FNN). We can see that the observation is  
 398 consistent with Figure 3a, where the alignment along eigen-directions is significant.

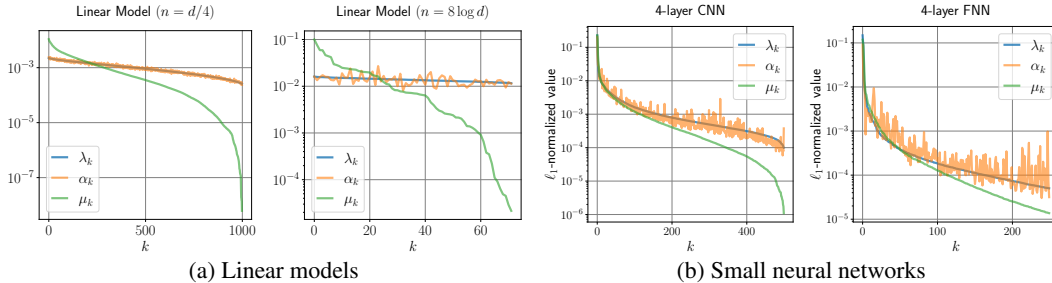


Figure 3: How the components of noise energy in *eigen-directions*  $\{\alpha_k\}_k$  are proportional to the corresponding curvatures  $\{\lambda_k\}_k$ .  $\alpha_k/\lambda_k$  can reflect the directional alignment along the eigen-directions of the local landscape. The eigenvalues of  $\Sigma/2\mathcal{L}$  are also plotted as comparison. (a) Linear models on Gaussian data in the regimes with limited data, where we fix  $d = 10^3$  and change  $n$  accordingly ( $n = d/4$ ,  $n = 8 \log d$ ). (b) 4-layer CNN and 4-layer FNN on CIFAR-10 dataset. For more experimental details, we refer to Appendix F.

### 399 C.3 Escaping Direction

400 Figure 4 presents numerical comparisons of the escaping directions between SGD and GD. It is  
 401 evident that  $D_1(t)$  exponentially decreases to zero for GD, indicating that GD escapes along the  
 402 sharpest direction. In contrast, for SGD,  $D_1(t)$  remains significantly large, indicating that SGD  
 403 retains a substantial component along the flat directions during the escape process. Furthermore,  
 404 the value of  $D_1(t)$  positively correlates with  $\text{srk}(G^2)$ , as predicted by our Theorem 5.2. These  
 405 observations provide empirical confirmation of our theoretical predictions.

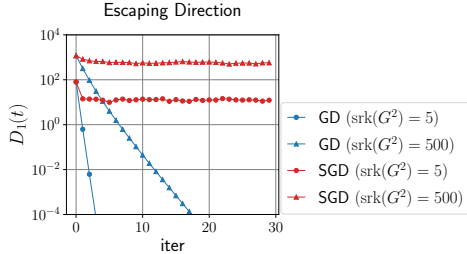


Figure 4: Comparison of escape directions between SGD and GD. The problem is linear regression and both SGD and GD are initialized near the global minimum by  $\mathbf{w}(0) \sim \mathcal{N}(\mathbf{w}^*, e^{-10} I_d/d)$ . To ensure escape, we choose  $\eta = 1.2/\|\mathbf{G}\|_F$  and  $\eta = 4/(\lambda_1 + \lambda_2)$  for SGD and GD, respectively. Please refer to Appendix F for more experimental details.

## 406 D Larger-scale Experiments for Deep Neural Networks

407 We have already provided small-scale experiments to confirm our theoretical findings. We now turn  
 408 to justify the practical relevance by examining the classification of CIFAR-10 dataset (Krizhevsky  
 409 and Hinton, 2009) with practical VGG nets (Simonyan and Zisserman, 2015) and ResNets (He et al.,  
 410 2016). Note that larger-scale experiments on average alignment have been previously presented in Wu  
 411 et al. (2022). Thus, our focus here is on investigate the directional alignment and escape direction of  
 412 SGD. We refer to Appendix F for experimental details.

413 **The directional alignment along eigen-directions.** Figure 5 presents the directional align-  
 414 ments of SGD noise for ResNet-38 and VGG-13. The alignment is examined along the eigen-  
 415 directions of the local landscape. The three quantities:  $\lambda_k$ ,  $\alpha_k$ , and  $\mu_k$  under  $\ell_1$  normalization  
 416 (i.e.,  $\lambda_k/\|\boldsymbol{\lambda}\|_1$ ,  $\alpha_k/\|\boldsymbol{\alpha}\|_1$ ,  $\mu_k/\|\boldsymbol{\mu}\|_1$ ) are plotted. Here,  $\lambda_k$  and  $\alpha_k$  represent the curvature and the  
 417 component of noise energy along the  $k$ -th eigen-direction, respectively.  $\mu_k$  corresponds to the  $k$ -th  
 418 eigenvalue of the noise covariance matrix, which is included for comparison. One can see that  
 419 the alignment between  $\alpha_k$  and  $\lambda_k$  still exists for ResNet-38 and VGG-13, but the ratio between  
 420 them becomes significantly larger. As a comparison, we refer to Figure 3b, where the ratio is well-  
 421 controlled for small-scale networks trained for classifying the same dataset. We hypothesize that this  
 422 observation is consistent with our theoretical results in Section 4: one-sided bounds require much  
 423 less samples.

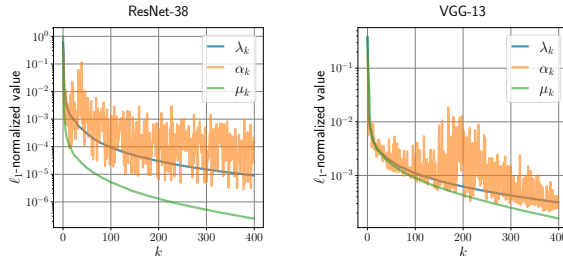


Figure 5: Three distributions ( $\{\lambda_k\}_k$ ,  $\{\alpha_k\}_k$ , and  $\{\mu_k\}_k$ ) for larger-scale neural networks, which reflect the directional alignment along the eigen directions of the local landscape.

424 **The escape direction of SGD.** For large models, it is computationally prohibitive to compute the  
 425 quantity  $D_k(t)$  since it needs to compute the whole spectrum. Thus, we consider to measure the  
 426 component along different directions without reweighting. Let  $\boldsymbol{\theta}^*$  be the minimum of interest  
 427 and  $\boldsymbol{\theta}(t)$  be SGD/GD solution at step  $t$ . Define  $p_k(t) = \langle \boldsymbol{\theta}(t) - \boldsymbol{\theta}^*, \mathbf{u}_1 \rangle$  for  $k = 1$  and  $p_k(t) =$   
 428  $(\sum_{i=1}^k \langle \boldsymbol{\theta}(t) - \boldsymbol{\theta}^*, \mathbf{u}_i \rangle^2)^{1/2}$  for  $k > 1$ ;  $r_k(t) = (\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}^*\|^2 - p_k^2(t))^{1/2}$ . Notably,  $p_k(t)$  and  $r_k(t)$   
 429 represent the component along sharp and flat directions, respectively.

430 In Figure 6, we plot  $(p_k(t), r_k(t))$  for VGG-19 and ResNet-110, where we examine various  $k$  values.  
 431 The plots clearly demonstrate that the escape direction of SGD exhibits significant components along  
 432 the flat directions. On the other hand, GD tends to escape along much sharper directions. These  
 433 empirical findings align well with our theoretical findings in Section 5.

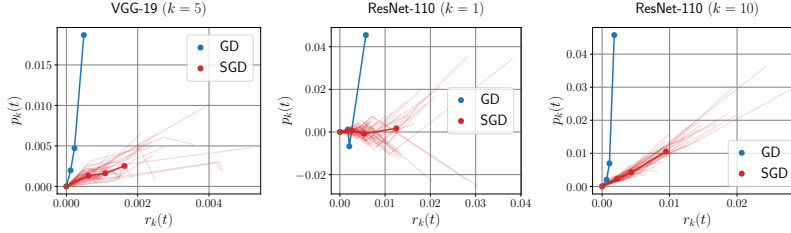


Figure 6: The red curves are 50 escaping trajectories of SGD and their average; the blue curves corresponding to GD. The sharp minimum  $\theta^*$  is found by SGD. Then, we run SGD and GD starting from  $\theta^*$  and the learning rates are tuned to ensure escaping.

## 434 E Conclusion and Future Work

435 In this paper, we present a comprehensive investigation of the geometry of SGD noise, demonstrating  
 436 both average and directional alignment between the noise and local geometry. We substantiate  
 437 these claims through both theoretical analyses and empirical evidence. Furthermore, we explore the  
 438 implications of these findings by analyzing the escape direction of SGD and its role in enhancing the  
 439 implicit bias toward flatter minima through cyclical learning rate.

440 Understanding the noise geometry is crucial for comprehending many aspects of stochastic optimiza-  
 441 tion, including but not limited to convergence rates, generalization capabilities, and dynamic behavior.  
 442 We offer an illustrative example through analyzing the escape direction of SGD. Another particularly  
 443 relevant application of our noise geometry framework lies in deciphering the Edge of Stability (EoS)  
 444 and the associated unstable convergence phenomena, as elaborated below.

445 • Studies (Cohen et al., 2020; Wu et al., 2018) showed that in training neural networks, GD typically  
 446 occurs in a EoS phase, where the the stability condition is violated. During EoS phase, GD  
 447 repeatedly slides into sharp regions and then, escapes from there. Due to the fact that GD escapes  
 448 along the sharpest direction (as stated in our Proposition 5.3), GD in the EoS phase will keep  
 449 oscillating along the sharpest directions and decreasing the loss along other flat directions. Thus,  
 450 EoS facilitates the unstable convergence of GD (Ahn et al., 2022). Similar EoS-related phenomena  
 451 and unstable convergence patterns are also observed in SGD (Lee and Jang, 2022). However, to  
 452 fully characterize the EoS phase in the context of SGD, it is imperative to understand the underlying  
 453 noise structure. Specifically, one must elucidate the mechanism by which noise compels SGD to  
 454 move away from sharp minima.

455 • In addition, our finding can potentially be used to explain why the training curve of SGD can be  
 456 more stable than that of GD—A very counter-intuitive phenomenon. As shown in Fig. 2 of Geiping  
 457 et al. (2021), GD training often encounters sudden large loss spikes and in contrast, SGD training  
 458 does not have this issue (although there are small loss fluctuations), implying that minibatch noise  
 459 can stabilize the training to some extent. This can potentially be explained by our theory as  
 460 follows. For both SGD and GD, the unstable dynamics is inevitable in training neural networks  
 461 due to progressive sharpening, i.e., entering the EoS phase. During the EoS phase, GD escapes  
 462 along the sharpest direction, leading to a sudden large loss spike if the curvature along the sharpest  
 463 direction becomes extremely large. In contrast, for SGD, the escape happens along much flatter  
 464 directions, for which it is unlikely to trigger a large loss spike.

## 465 F Experimental Setups

466 In this section, we provide the experiment details for directional alignment experiments (in Figure 3  
 467 and Figure 5) and escaping experiments (in Figure 4 and Figure 6).

468 **Small-scale experiments** (Figure 3 and 4).

469 • In Figure 3, we conduct experiments on linear regression and a 4-layer linear network:  $d \rightarrow$   
 470  $m \rightarrow m \rightarrow m \rightarrow 1$  with  $m = 50$ . The inputs  $\{x_i\}_{i=1}^n$  are drawn from  $\mathcal{N}(\mathbf{0}, I_d)$ . In the first  
 471 three experiments, we fix  $d = 10^3$  and change  $n$  accordingly ( $n = 4d^2, n = d, n = d/4$ ).

472 For the last experiment, we set  $d = 10^4$  and  $n = \log d$ . Regarding the parameter  $\theta$ , it is  
 473 drawn from  $\mathcal{N}(\mathbf{0}, I_p)$ .

- 474 • In Figure 4, we conduct escaping experiments on linear regression with  $\mathbf{w}^* = \mathbf{0}$ . Both SGD  
 475 and GD are initialized near the global minimum by  $\mathbf{w}(0) \sim \mathcal{N}(\mathbf{0}, e^{-10}I_d/d)$ . To ensure  
 476 escaping, we choose  $\eta = 1.2/\|G\|_F$  and  $\eta = 4/(\lambda_1 + \lambda_2)$  for SGD and GD, respectively.  
 477 We fix  $n = 10^5$  and  $d = 10^3$ , and the inputs  $\{\mathbf{x}_i\}_{i=1}^n$  are drawn from  $\mathcal{N}(\mathbf{0}, \text{diag}(\boldsymbol{\lambda})/d)$ ,  
 478 where  $\boldsymbol{\lambda} \in \mathbb{R}^d$  and  $\lambda_1 \geq \lambda_2 = \dots = \lambda_d \geq 0$ . Moreover, we set  $\lambda_1 = 1$  change  $\lambda_2$   
 479 accordingly to obtain different  $\text{srk}(G^2)$ .

480 **Larger-scale experiments** (Figure 5 and 6).

- 481 • **Dataset.** For the experiments in Figure 5 and 6, we use the CIFAR-10 dataset with label=0, 1  
 482 and the full CIFAR-10 dataset to train our models, respectively.

- 483 • **Models.** We conduct experiments on large-scale models: 4-layer CNN ( $p = 43, 072$ ),  
 484 4-layer FNN ( $p = 219, 200$ ), ResNet-38 ( $p = 558, 222$ ), VGG-13 ( $p = 605, 458$ ), ResNet-  
 485 110 ( $p = 1, 720, 138$ ), and VGG-19 ( $p = 20, 091, 338$ ).

486 Specifically, we use standard ResNets (He et al., 2016) and VGG nets (Simonyan and  
 487 Zisserman, 2015) without batch normalization. For ResNets, we follow Zhang et al. (2019)  
 488 to use the fixup initialization in order to ensure that the model can be trained without batch  
 489 normalization. Moreover, the architecture of 4-layer CNN is  $\text{Conv}(3, 6, 5) \rightarrow \text{ReLU} \rightarrow$   
 490  $\text{MPool}(2, 2) \rightarrow \text{Conv}(6, 16, 5) \rightarrow \text{ReLU} \rightarrow \text{MPool}(2, 2) \rightarrow \text{Linear}(400, 100) \rightarrow \text{ReLU} \rightarrow$   
 491  $\text{Linear}(100, 2)$ . and the 4-layer FNN is a ReLU-activated fully-connected network with  
 492 the architecture:  $784 \rightarrow 256 \rightarrow 64 \rightarrow 32 \rightarrow 2$ .

- 493 • **Training.** All explicit regularizations (including weight decay, dropout, data augmentation,  
 494 batch normalization, learning rate decay) are removed, and a simple constant-LR SGD is  
 495 used to train our models. Specifically, all these models are trained by SGD with learning  
 496 rate  $\eta = 0.1$  and batch size  $B = 32$  until the training loss becomes smaller than  $10^{-4}$ .

497 **Efficient computations** of the top- $k$  eigen-decomposition of  $G$  and  $\Sigma$ . We utilize the functions `eighs`  
 498 and `LinearOperator` in `scipy.sparse.linalg` to calculate top- $k$  eigenvalues and eigenvectors  
 499 of  $G$  and  $\Sigma$ , and the key step is to efficiently calculate  $G\mathbf{v}$  and  $\Sigma\mathbf{v}$  for any given  $\mathbf{v} \in \mathbb{R}^p$ .

- 500 • For small-scale experiments, they can be calculated directly.
- 501 • For the large-scale models, we need further approximations since the computation complex-  
 502 ity  $\mathcal{O}(np)$  is prohibitive in this case. To illustrate our method, we will use  $G\mathbf{v}$  as an example  
 503 and apply a similar approach to  $\Sigma\mathbf{v}$ . Notice that the formulation  $G\mathbf{v} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{v}) \mathbf{x}_i$   
 504 are all in the form of sample average, which allows us to perform Monte-Carlo approxi-  
 505 mation. Specifically, we randomly choose  $b$  samples  $\{\mathbf{x}_{i_j}\}_{j=1}^b$  from  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and use  
 506  $\frac{1}{b} \sum_{j=1}^b (\mathbf{x}_{i_j}^\top \mathbf{v}) \mathbf{x}_{i_j}$  estimate  $G\mathbf{v}$ , with the computation complexity  $\mathcal{O}(bp)$ . For the experi-  
 507 ments on CIFAR-10, we test  $b$ 's with different values and find that  $b = 2k$  is sufficient to  
 508 obtain a reliable approximation of the top- $k$  eigenvalues and eigenvectors. Hence, for all  
 509 large-scale experiments in this paper, we use  $b = 2k$  to speed up the computation of the  
 510 top- $k$  eigenvalues and eigenvectors.

## 511 G Proofs in Section 3: Average alignment

### 512 G.1 Proof of Theorem 3.1 (a)

513 For clarity, in a slightly different order from the main text, we first prove for the linear model  
 514 (Example) and then for the OLM (Theorem 3.1). This is also convenient for us to compare the  
 515 difference between the proof for the two-layer neural network (Theorem 3.4) and the proof for the  
 516 linear model.

517 **Step I.** *Proof for linear models.*

518 For the linear model, i.e.,  $\boldsymbol{\theta} = \boldsymbol{w}$  and  $F(\boldsymbol{w}) = \boldsymbol{w}$  in OLMs, we have

$$\begin{aligned}
\mu(\boldsymbol{w}) &= \frac{\text{Tr}(\Sigma(\boldsymbol{w})G(\boldsymbol{w}))}{2\mathcal{L}(\boldsymbol{w})\|G(\boldsymbol{w})\|_F^2} \\
&= \frac{\text{Tr}\left(\left(\frac{1}{n}\sum_{j=1}^n \boldsymbol{x}_j \boldsymbol{x}_j^\top\right)\left(\frac{1}{n}\sum_{i=1}^n (F(\boldsymbol{\theta})^\top \boldsymbol{x}_i)^2 (\nabla F(\boldsymbol{\theta})^\top \boldsymbol{x}_i)(\nabla F(\boldsymbol{\theta})^\top \boldsymbol{x}_i)^\top\right)\right)}{\left(\frac{1}{n}\sum_{i=1}^n (F(\boldsymbol{\theta})^\top \boldsymbol{x}_i)^2\right)\left(\frac{1}{n^2}\sum_{i=1}^n \sum_{j=1}^n (\boldsymbol{x}_i^\top \nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top \boldsymbol{x}_j)^2\right)} \\
&= \frac{\frac{1}{n^2}\sum_{i=1}^n \sum_{j=1}^n (\boldsymbol{w}^\top \boldsymbol{x}_i)^2 (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2}{\left(\frac{1}{n}\sum_{i=1}^n (\boldsymbol{w}^\top \boldsymbol{x}_i)^2\right)\left(\frac{1}{n^2}\sum_{i=1}^n \sum_{j=1}^n (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2\right)} \geq \frac{\left(\frac{1}{n}\sum_{i=1}^n (\boldsymbol{w}^\top \boldsymbol{x}_i)^2\right)\left(\min_{i \in [n]} \frac{1}{n}\sum_{j=1}^n (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2\right)}{\left(\frac{1}{n}\sum_{i=1}^n (\boldsymbol{w}^\top \boldsymbol{x}_i)^2\right)\left(\frac{1}{n^2}\sum_{i=1}^n \sum_{j=1}^n (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2\right)} \quad (1) \\
&= \frac{\min_{i \in [n]} \frac{1}{n}\sum_{j=1}^n (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2}{\max_{i \in [n]} \frac{1}{n}\sum_{j=1}^n (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2} \geq \frac{\min_{i \in [n]} \|\boldsymbol{x}_i\|^4 + (n-1)\min_{i \in [n]} \frac{1}{n-1}\sum_{j \neq i} (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2}{\max_{i \in [n]} \|\boldsymbol{x}_i\|^4 + (n-1)\max_{i \in [n]} \frac{1}{n-1}\sum_{j \neq i} (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2}.
\end{aligned}$$

519 Then we only need to estimate  $\|\boldsymbol{x}_i\|^4$  and  $\frac{1}{n-1}\sum_{j \neq i} (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2$  for each  $i \in [n]$ , respectively.

520 Step I (i). Estimation of  $\|\boldsymbol{x}_i\|^4$ .

521 Let  $\boldsymbol{y}_i = S^{1/2}\boldsymbol{x}_i$ , then  $\|\boldsymbol{x}_i\|^2 = \boldsymbol{y}_i^\top S \boldsymbol{y}_i$  and  $\boldsymbol{y}_1, \dots, \boldsymbol{y}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I_d)$ .

522 For a fix  $i \in [n]$ , by Lemma J.2, there exists an absolute constant  $C_1 > 0$  such that for any  $\epsilon \in (0, 1)$ ,  
523 we have

$$\mathbb{P}\left(\left|\boldsymbol{y}_i^\top S \boldsymbol{y}_i - \text{Tr}(S)\right| \geq \epsilon \text{Tr}(S)\right) \leq 2 \exp\left(-C_1 \min\left\{\frac{\epsilon^2 \text{Tr}^2(S)}{\|S\|_F^2}, \frac{\epsilon \text{Tr}(S)}{\|S\|_2}\right\}\right).$$

524 Noticing that  $\text{Tr}(S)\|S\|_2 = \lambda_1 \sum_i \lambda_i \geq \sum_i \lambda_i^2 = \|S\|_F$ , we thus have

$$\frac{\text{Tr}^2(S)}{\|S\|_F^2} \geq \frac{\text{Tr}(S)}{\|S\|_2} = \text{srk}(S).$$

525 Therefore,

$$\mathbb{P}\left(\left|\boldsymbol{y}_i^\top S \boldsymbol{y}_i - \text{Tr}(S)\right| \geq \epsilon \text{Tr}(S)\right) \leq 2 \exp\left(-C_1 \frac{\text{Tr}(S)}{\|S\|_2} \min\{\epsilon, \epsilon^2\}\right) = 2 \exp(-C_1 \epsilon^2 \text{srk}(S)).$$

526 Applying a union bound over all  $i \in [n]$ , we have

$$\mathbb{P}\left(\left|\|\boldsymbol{x}_i\|^2 - \text{Tr}(S)\right| \geq \epsilon \text{Tr}(S), \forall i \in [n]\right) \leq 2n \exp(-C_1 \epsilon^2 \text{srk}(S)).$$

527 In the other word, for any  $\epsilon, \delta \in (0, 1)$ , if  $\text{srk}(S) \gtrsim \log(n)/\epsilon^2$ , then *w.p.* at least  $1 - \delta/3$ , we have

$$(1 - \epsilon)^2 \leq \frac{\|\boldsymbol{x}_i\|_2^4}{\text{Tr}^2(S)} \leq (1 + \epsilon)^2, \forall i \in [n].$$

528 Step I (ii). Estimation of  $\frac{1}{n-1}\sum_{j \neq i} (\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2$ .

529 First, we fix  $i \in [n]$ . Notice that  $(\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2$  ( $j \neq i$ ) are not independent, so we need estimate by some  
530 decoupling tricks.

531 We denote  $\boldsymbol{y}_i := S^{-1/2}\boldsymbol{x}_i$ , then  $\boldsymbol{y}_1, \dots, \boldsymbol{y}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I_d)$  and  $(\boldsymbol{x}_i^\top \boldsymbol{x}_j)^2 = (\boldsymbol{y}_i^\top S \boldsymbol{y}_j)^2$ .

532 For any fixed  $\mathbf{v} \in \mathbb{S}^{d-1}$ , by Lemma J.1, for any  $\epsilon \in (0, 1)$ , we have

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{v}^\top \mathbf{y}_j)^2 - 1 \right| \geq \epsilon \right) \\ & \leq \mathbb{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{v}^\top \mathbf{y}_j)^2 - 1 \right| \geq \epsilon \right) \leq 2 \exp(-C_2(n-1)\epsilon^2), \end{aligned}$$

533 where  $C_2 > 0$  is an absolute constant, independent of  $\mathbf{v}$  and  $\epsilon$ .

534 Then we have

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{x}_i^\top \mathbf{x}_j)^2 - \mathbf{x}_i^\top S \mathbf{x}_i \right| \geq \epsilon \mathbf{x}_i^\top S \mathbf{x}_i \right) \\ & = \mathbb{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{y}_i^\top S \mathbf{y}_j)^2 - \|S \mathbf{y}_i\|_2^2 \right| \geq \epsilon \|S \mathbf{y}_i\|_2^2 \right) \\ & \stackrel{\mathbf{z}_i := S \mathbf{y}_i / \|S \mathbf{y}_i\|_2}{=} \mathbb{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{z}_i^\top \mathbf{y}_j)^2 - 1 \right| \geq \epsilon \right) \\ & = \mathbb{E} \left[ \mathbb{I} \left\{ \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{z}_i^\top \mathbf{y}_j)^2 - 1 \right| \geq 1 \right\} \right] \\ & = \mathbb{E}_{\mathbf{z}_i} \left[ \mathbb{E} \left[ \mathbb{I} \left\{ \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{z}_i^\top \mathbf{y}_j)^2 - 1 \right| \geq 1 \right\} \middle| \mathbf{z}_i \right] \right] \\ & \leq \mathbb{E}_{\mathbf{z}_i} [2 \exp(-C_2(n-1)\epsilon^2)] = 2 \exp(-C_2(n-1)\epsilon^2). \end{aligned}$$

535 Applying a union bound over all  $i \in [n]$ , we have

$$\mathbb{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} (\mathbf{x}_i^\top \mathbf{x}_j)^2 - \mathbf{x}_i^\top S \mathbf{x}_i \right| \geq \epsilon \mathbf{x}_i^\top S \mathbf{x}_i, \forall i \in [n] \right) \leq 2n \exp(-C_2(n-1)\epsilon^2).$$

536 In the other word, for any  $\epsilon, \delta \in (0, 1)$ , if  $n/\log(n/\delta) \gtrsim 1/\epsilon^2$ , then *w.p.* at least  $1 - \delta/3$ , we have

$$1 - \epsilon \leq \frac{\frac{1}{n-1} \sum_{j \neq i} (\mathbf{x}_i^\top \mathbf{x}_j)^2}{\mathbf{x}_i^\top S \mathbf{x}_i} \leq 1 + \epsilon, \forall i \in [n].$$

537 Step I (iii). Estimation of  $\mathbf{x}_i^\top S \mathbf{x}_i$ .

538 Let  $\mathbf{y}_i = S^{1/2} \mathbf{x}_i$ , then  $\mathbf{x}_i^\top S \mathbf{x}_i = \mathbf{y}_i^\top S^2 \mathbf{y}_i$  and  $\mathbf{y}_1, \dots, \mathbf{y}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I_d)$ .

539 In the same way as Step I(i), we obtain that: for any  $\epsilon, \delta \in (0, 1)$ , if  $\text{srk}(S^2) \gtrsim \log(n)/\epsilon^2$ , then  
540 *w.p.* at least  $1 - \delta/3$ , we have

$$1 - \epsilon \leq \frac{\mathbf{x}_i^\top S \mathbf{x}_i}{\text{Tr}(S^2)} \leq 1 + \epsilon, \forall i \in [n].$$

541 Combining our results in Step I (i), Step I (ii), and Step I (iii), we obtain the result for Linear Model:  
542 for any  $\epsilon, \delta \in (0, 1)$ , if  $n/\log(n/\delta) \gtrsim 1/\epsilon^2$  and  $\min\{\text{srk}(S), \text{srk}(S^2)\} \gtrsim \log(n)/\epsilon^2$ , then *w.p.* at  
543 least  $1 - \delta/3 - \delta/3 - \delta/3 = 1 - \delta$ , we have

$$\mu(\mathbf{w}) \geq \frac{(1 - \epsilon)^2 \text{Tr}^2(S) + (n-1)(1 - \epsilon) \min_{i \in [n]} \mathbf{x}_i^\top S \mathbf{x}_i}{(1 + \epsilon)^2 \text{Tr}^2(S) + (n-1)(1 + \epsilon) \max_{i \in [n]} \mathbf{x}_i^\top S \mathbf{x}_i}$$



$$\geq \frac{(1-\epsilon)^2 \text{Tr}^2(S) + (n-1)(1-\epsilon)^2 \text{Tr}(S^2)}{(1+\epsilon)^2 \text{Tr}^2(S) + (n-1)(1+\epsilon)^2 \text{Tr}(S^2)} = \frac{(1-\epsilon)^2}{(1+\epsilon)^2}.$$

544 From the arbitrary of  $\mathbf{w}$ , we have  $\inf_{\mathbf{w} \in \mathbb{R}^d} \mu(\mathbf{w}) \geq \frac{(1-\epsilon)^2}{(1+\epsilon)^2}$ .

545 **Step II. Proof for OLMs.**

$$\begin{aligned} \mu(\boldsymbol{\theta}) &= \frac{\text{Tr}(\Sigma(\boldsymbol{\theta})G(\boldsymbol{\theta}))}{2\mathcal{L}(\boldsymbol{\theta})\|G(\boldsymbol{\theta})\|_{\text{F}}^2} \\ &= \frac{\text{Tr}\left(\left(\frac{1}{n}\sum_{j=1}^n(\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)(\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^\top\right)\left(\frac{1}{n}\sum_{i=1}^n(F(\boldsymbol{\theta})^\top \mathbf{x}_i)^2(\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i)(\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i)^\top\right)\right)}{\left(\frac{1}{n}\sum_{i=1}^n(F(\boldsymbol{\theta})^\top \mathbf{x}_i)^2\right)\left(\frac{1}{n^2}\sum_{i=1}^n\sum_{j=1}^n(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2\right)} \\ &= \frac{\frac{1}{n^2}\sum_{i=1}^n\sum_{j=1}^n(F(\boldsymbol{\theta})^\top \mathbf{x}_i)^2(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2}{\left(\frac{1}{n}\sum_{i=1}^n(F(\boldsymbol{\theta})^\top \mathbf{x}_i)^2\right)\left(\frac{1}{n^2}\sum_{i=1}^n\sum_{j=1}^n(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2\right)} \\ &\geq \frac{\left(\frac{1}{n}\sum_{i=1}^n(F(\boldsymbol{\theta})^\top \mathbf{x}_i)^2\right)\left(\min_{i \in [n]} \frac{1}{n}\sum_{j=1}^n(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2\right)}{\left(\frac{1}{n}\sum_{i=1}^n(F(\boldsymbol{\theta})^\top \mathbf{x}_i)^2\right)\left(\frac{1}{n^2}\sum_{i=1}^n\sum_{j=1}^n(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2\right)} = \frac{\min_{i \in [n]} \frac{1}{n}\sum_{j=1}^n(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2}{\max_{i \in [n]} \frac{1}{n}\sum_{j=1}^n(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2} \\ &\geq \frac{\min_{i \in [n]} \|\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i\|^4 + (n-1)\min_{i \in [n]} \frac{1}{n-1}\sum_{j \neq i}(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2}{\max_{i \in [n]} \|\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i\|^4 + (n-1)\max_{i \in [n]} \frac{1}{n-1}\sum_{j \neq i}(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2}. \end{aligned} \tag{2}$$

546 We can still prove the theorem by the similar way as Step I.

547 By replacing  $\mathbf{x}_i$  and  $\mathbf{x}_j$  ( $j \neq i$ ) in Step I (i) with  $\nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i$  and  $\mathbf{x}_j$  ( $j \neq i$ ), respectively, in  
548 the similar way as Step I (i), we can obtain: for any  $\epsilon, \delta \in (0, 1)$ , if  $n/\log(n/\delta) \gtrsim 1/\epsilon^2$ , then *w.p.* at  
549 least  $1 - \delta$ , we have

$$1 - \epsilon \leq \frac{\frac{1}{n-1}\sum_{j \neq i}(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2}{\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top S \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i} \leq 1 + \epsilon, \quad \forall i \in [n];$$

550 Combining the estimation above with Step I (ii) and Step I (iii), we obtain that: for any  $\epsilon, \delta \in (0, 1)$ ,  
551 if  $n/\log(n/\delta) \gtrsim 1/\epsilon^2$  and  $\text{srk}(S^2) \gtrsim \log(n)/\epsilon^2$ , then *w.p.* at least  $1 - \delta$ , we have

$$\begin{aligned} 1 - \epsilon &\leq \frac{\frac{1}{n-1}\sum_{j \neq i}(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2}{\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top S \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i} \leq 1 + \epsilon, \quad \forall i \in [n]; \\ (1 - \epsilon)^2 &\leq \frac{\|\mathbf{x}_i\|_2^4}{\text{Tr}^2(S)} \leq (1 + \epsilon)^2, \quad \forall i \in [n]; \\ 1 - \epsilon &\leq \frac{\mathbf{x}_i^\top S \mathbf{x}_i}{\text{Tr}(S^2)} \leq 1 + \epsilon, \quad \forall i \in [n]. \end{aligned}$$

552 These inequalities imply that:

$$\begin{aligned} \mu(\boldsymbol{\theta}) &\geq \frac{\min_{i \in [n]} \lambda_{\min}^2(\nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top) \|\mathbf{x}_i\|_2^4 + (n-1)\min_{i \in [n]} \frac{1}{n-1}\sum_{j \neq i}(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2}{\max_{i \in [n]} \lambda_{\min}^2(\nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top) \|\mathbf{x}_i\|_2^4 + (n-1)\max_{i \in [n]} \frac{1}{n-1}\sum_{j \neq i}(\mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_j)^2} \\ &\geq \frac{(1 - \epsilon)^2 \min_{i \in [n]} \lambda_{\min}^2(\nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top) \text{Tr}^2(S) + (n-1)(1 - \epsilon) \min_{i \in [n]} \mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top S \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i}{(1 - \epsilon)^2 \max_{i \in [n]} \lambda_{\max}^2(\nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top) \text{Tr}^2(S) + (n-1)(1 + \epsilon) \max_{i \in [n]} \mathbf{x}_i^\top \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top S \nabla F(\boldsymbol{\theta})\nabla F(\boldsymbol{\theta})^\top \mathbf{x}_i} \end{aligned}$$

$$\begin{aligned}
& \frac{(1-\epsilon)^2 \min_{i \in [n]} \lambda_{\min}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \text{Tr}^2(S) + (n-1)(1-\epsilon) \lambda_{\min}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \min_{i \in [n]} \mathbf{x}_i^\top S \mathbf{x}_i}{(1+\epsilon)^2 \max_{i \in [n]} \lambda_{\max}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \text{Tr}^2(S) + (n-1)(1+\epsilon) \lambda_{\max}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \max_{i \in [n]} \mathbf{x}_i^\top S \mathbf{x}_i} \\
& \geq \frac{(1-\epsilon)^2 \min_{i \in [n]} \lambda_{\min}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \text{Tr}^2(S) + (n-1)(1-\epsilon)^2 \lambda_{\min}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \text{Tr}(S^2)}{(1+\epsilon)^2 \max_{i \in [n]} \lambda_{\max}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \text{Tr}^2(S) + (n-1)(1+\epsilon)^2 \lambda_{\max}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top) \text{Tr}(S^2)} \\
& = \frac{(1-\epsilon)^2}{(1+\epsilon)^2 \text{cond}^2(\nabla F(\boldsymbol{\theta}) \nabla F(\boldsymbol{\theta})^\top)}.
\end{aligned}$$

553 Hence, we have proved Theorem 3.1.  $\square$

## 554 G.2 Proof of Theorem 3.1 (b)

555 This result is a direct corollary of Theorem 4.2, which is proved in Appendix H.

556 Under the same setting as Theorem 4.2, Theorem 4.2 gives us the uniform lower bound: there exists  
557 an absolute constant  $C > 0$  such that

$$\inf_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p} g(\boldsymbol{\theta}; \mathbf{v}) \geq C,$$

558 which means that for any  $\boldsymbol{\theta} \in \mathbb{R}^p$ ,  $\mathbf{v} \in \mathbb{S}^{p-1}$ , we have

$$\mathbf{v}^\top \Sigma(\boldsymbol{\theta}) \mathbf{v} \geq C \cdot 2\mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}.$$

559 Consider the orthogonal decomposition of  $G(\boldsymbol{\theta})$ :  $G(\boldsymbol{\theta}) = \sum_{k=1}^p \lambda_k \mathbf{u}_k \mathbf{u}_k^\top$ . Notice that

$$\begin{aligned}
\text{Tr}(\Sigma(\boldsymbol{\theta}) G(\boldsymbol{\theta})) &= \sum_{k=1}^p \lambda_k \mathbf{u}_k^\top \Sigma(\boldsymbol{\theta}) \mathbf{u}_k, \\
\|G(\boldsymbol{\theta})\|_F &= \text{Tr}(G(\boldsymbol{\theta}) G(\boldsymbol{\theta})) = \sum_{k=1}^p \lambda_k \mathbf{u}_k^\top G(\boldsymbol{\theta}) \mathbf{u}_k.
\end{aligned}$$

560 Then we obtain

$$\text{Tr}(\Sigma(\boldsymbol{\theta}) G(\boldsymbol{\theta})) \geq C \cdot 2\mathcal{L}(\boldsymbol{\theta}) \sum_{k=1}^p \lambda_k \mathbf{u}_k^\top G(\boldsymbol{\theta}) \mathbf{u}_k = C \cdot 2\mathcal{L}(\boldsymbol{\theta}) \|G(\boldsymbol{\theta})\|_F^2,$$

561 which means  $\mu(\boldsymbol{\theta}) \geq C$ . From the arbitrariness of  $\boldsymbol{\theta}$ , it holds that  $\inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \mu(\boldsymbol{\theta}) \geq C$ .  $\square$

## 562 G.3 Proof of Theorem 3.4

563 For two-layer neural networks with fixed output layer, the gradient is

$$\nabla f(\mathbf{x}_i; \boldsymbol{\theta}) = \left( a_1 \sigma'(\mathbf{b}_1^\top \mathbf{x}_i) \mathbf{x}_i^\top, \dots, a_m \sigma'(\mathbf{b}_m^\top \mathbf{x}_i) \mathbf{x}_i^\top \right)^\top \in \mathbb{R}^{md}.$$

564 For simplicity, denote  $\nabla f_i(\boldsymbol{\theta}) := \nabla f(\mathbf{x}_i; \boldsymbol{\theta})$ ,  $\mathbf{u}_i(\boldsymbol{\theta}) := f_i(\boldsymbol{\theta}) - f_i(\boldsymbol{\theta}^*)$ . Then we have:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n u_i^2(\boldsymbol{\theta}), \quad G(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}) \nabla f_i(\boldsymbol{\theta})^\top, \quad \Sigma(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n u_i^2(\boldsymbol{\theta}) \nabla f_i(\boldsymbol{\theta}) \nabla f_i(\boldsymbol{\theta})^\top.$$

$$\mu(\boldsymbol{\theta}) = \frac{\text{Tr} \left( \left( \frac{1}{n} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}) \nabla f_i(\boldsymbol{\theta})^\top \right) \left( \frac{1}{n} \sum_{i=1}^n u_i^2(\boldsymbol{\theta}) \nabla f_i(\boldsymbol{\theta}) \nabla f_i(\boldsymbol{\theta})^\top \right) \right)}{\left( \frac{1}{n} \sum_{i=1}^n u_i^2(\boldsymbol{\theta}) \right) \left( \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\nabla f_i(\boldsymbol{\theta})^\top \nabla f_j(\boldsymbol{\theta}))^2 \right)}$$

$$\begin{aligned}
& \frac{\frac{1}{n} \sum_{i=1}^n u_i^2(\boldsymbol{\theta}) \frac{1}{n} \sum_{j=1}^n (\nabla f_i(\boldsymbol{\theta})^\top \nabla f_j(\boldsymbol{\theta}))^2}{\left( \frac{1}{n} \sum_{i=1}^n u_i^2(\boldsymbol{\theta}) \right) \left( \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\nabla f_i(\boldsymbol{\theta})^\top \nabla f_j(\boldsymbol{\theta}))^2 \right)} \\
& \geq \frac{\min_{i \in [n]} \frac{1}{n} \sum_{j=1}^n (\nabla f_i(\boldsymbol{\theta})^\top \nabla f_j(\boldsymbol{\theta}))^2}{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\nabla f_i(\boldsymbol{\theta})^\top \nabla f_j(\boldsymbol{\theta}))^2} \geq \frac{\min_{i \in [n]} \frac{1}{n} \sum_{j=1}^n (\alpha^2 m \mathbf{x}_i^\top \mathbf{x}_j)^2}{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\beta^2 m \mathbf{x}_i^\top \mathbf{x}_j)^2} = \frac{\alpha^2 \min_{i \in [n]} \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_i^\top \mathbf{x}_j)^2}{\beta^2 \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i^\top \mathbf{x}_j)^2}.
\end{aligned}$$

565 Notice that the last term  $\frac{\min_{i \in [n]} \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_i^\top \mathbf{x}_j)^2}{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i^\top \mathbf{x}_j)^2}$  is independent of  $\boldsymbol{\theta}$  and the same as (1) for the linear  
566 model. Then repeating the same proof of Linear Model, the result of this theorem differs from Linear  
567 Model by only the factor  $\alpha^2/\beta^2$ . In other words, under the same condition with Linear Model, w.p. at  
568 least  $1 - \delta$ , we have

$$\inf_{\boldsymbol{\theta} \in \mathbb{R}^{md}} \mu(\boldsymbol{\theta}) \geq \frac{\alpha^2 (1 - \epsilon)^2}{\beta^2 (1 + \epsilon)^2}.$$

569

□

## 570 H Proofs in Section 4: Directional Alignment

571 For the OLM  $f(\mathbf{x}; \boldsymbol{\theta}) = F(\boldsymbol{\theta})^\top \mathbf{x}$ , let  $\mathbf{r}(\boldsymbol{\theta}) = F(\boldsymbol{\theta}) - F(\boldsymbol{\theta}^*)$ . Then, we have

$$\begin{aligned}
\hat{G}(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \nabla F^\top(\boldsymbol{\theta}) \mathbf{x}_i \mathbf{x}_i^\top \nabla F(\boldsymbol{\theta}) \\
\hat{\mathcal{L}}(\boldsymbol{\theta}) &= \frac{1}{2n} \sum_{i=1}^n (\mathbf{u}^\top(\boldsymbol{\theta}) \mathbf{x}_i)^2 \\
\hat{\Sigma}(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n (\mathbf{r}^\top(\boldsymbol{\theta}) \mathbf{x}_i)^2 \nabla F^\top(\boldsymbol{\theta}) \mathbf{x}_i \mathbf{x}_i^\top \nabla F(\boldsymbol{\theta}),
\end{aligned} \tag{3}$$

572 and for the population case:

$$\begin{aligned}
G(\boldsymbol{\theta}) &= \mathbb{E} \left[ \nabla F^\top(\boldsymbol{\theta}) \mathbf{x} \mathbf{x}^\top \nabla F(\boldsymbol{\theta}) \right] = \nabla F^\top(\boldsymbol{\theta}) S \nabla F(\boldsymbol{\theta}) \\
\mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E} \left[ (\mathbf{r}^\top(\boldsymbol{\theta}) \mathbf{x})^2 \right] = \frac{1}{2} \mathbf{r}(\boldsymbol{\theta})^\top S \mathbf{r}(\boldsymbol{\theta}) \\
\Sigma(\boldsymbol{\theta}) &= \mathbb{E} \left[ (\mathbf{r}^\top(\boldsymbol{\theta}) \mathbf{x})^2 \nabla F^\top(\boldsymbol{\theta}) \mathbf{x} \mathbf{x}^\top \nabla F(\boldsymbol{\theta}) \right]
\end{aligned}$$

**Lemma H.1** (Proposition 2.3 in (Wu et al., 2022)). *Let the data distribution be  $\mathcal{N}(\mathbf{0}, S)$ . Then we have*

$$\Sigma(\boldsymbol{\theta}) = \nabla \mathcal{L}(\boldsymbol{\theta}) \nabla \mathcal{L}(\boldsymbol{\theta})^\top + 2\mathcal{L}(\boldsymbol{\theta}) G(\boldsymbol{\theta}).$$

573 **Lemma H.2.** *Under the same conditions in Lemma H.1, if  $\mathbf{u}(\boldsymbol{\theta}) \neq \mathbf{0}$  and  $\nabla F(\boldsymbol{\theta}) \mathbf{v} \neq \mathbf{0}$ , then we*  
574 *have:*

$$(\nabla \mathcal{L}(\boldsymbol{\theta})^\top \mathbf{v})^2 \leq 2\mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}.$$

575 *Proof.* Noticing that  $\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbf{r}(\boldsymbol{\theta})^\top S \mathbf{r}(\boldsymbol{\theta})$ , we have  $\nabla \mathcal{L}(\boldsymbol{\theta}) = \nabla F(\boldsymbol{\theta})^\top S \mathbf{u}(\boldsymbol{\theta})$ . Hence,

$$\begin{aligned}
& (\nabla \mathcal{L}(\boldsymbol{\theta})^\top \mathbf{v})^2 = \mathbf{v}^\top \nabla F(\boldsymbol{\theta})^\top S \mathbf{r}(\boldsymbol{\theta}) \mathbf{r}(\boldsymbol{\theta})^\top S \nabla F(\boldsymbol{\theta}) \mathbf{v} = \langle \nabla F(\boldsymbol{\theta}) \mathbf{v}, \mathbf{r}(\boldsymbol{\theta}) \rangle_S^2 \\
& \stackrel{\text{Lemma J.6}}{\leq} \|\nabla F(\boldsymbol{\theta}) \mathbf{v}\|_S^2 \|\mathbf{r}(\boldsymbol{\theta})\|_S^2 = 2\mathcal{L}(\boldsymbol{\theta}) \left( \mathbf{v} \nabla F(\boldsymbol{\theta})^\top S \nabla F(\boldsymbol{\theta}) \mathbf{v} \right) = 2\mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}.
\end{aligned}$$

576

□

577 **Lemma H.3.** Let  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . For any  $\epsilon, \delta \in (0, 1)$ , if we choose  $n \gtrsim$   
578  $(d + \log(1/\delta)) / \epsilon^2$ , then w.p. at least  $1 - \delta$ , we have:

$$\sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \mathbf{x}_i)^2 - 1 \right| \leq \epsilon.$$

579 *Proof.* By Lemma J.3 with  $K = \sqrt{C_1}$ , we know that: w.p. at least  $1 - 2 \exp(-u)$ , we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \mathbf{I}_d \right\| \leq C_2 \left( \sqrt{\frac{d+u}{n}} + \frac{d+u}{n} \right),$$

580 where  $C_2$  is an absolute positive constant. Equivalently, we can rewrite this conclusion. For any  
581  $\epsilon, \delta \in (0, 1)$ , if we choose  $n \gtrsim (d + \log(1/\delta)) / \epsilon^2$ , then w.p. at least  $1 - \delta$ , we have:

$$\sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \mathbf{x}_i)^2 - 1 \right| \leq \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \mathbf{I}_d \right\| \leq \epsilon.$$

582 □

583 **Lemma H.4** (Corollary 2 in (Cai et al., 2022)). Let  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . There exists absolute  
584 constants  $C_1, C_2, C_3 > 0$ , such that if  $n \geq C_3 d$ , then w.p. at least  $1 - \exp(-C_2 n)$ , we have

$$\inf_{\mathbf{u}, \mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{u})^2 (\mathbf{x}_i^\top \mathbf{v})^2 \geq C_1.$$

585 With the preparation of Lemma H.3 and Lemma H.4, now we give the proof of Theorem 4.2.

### 586 H.1 Proof of Theorem 4.2

587 Let  $\mathbf{y}_i = \mathbf{S}^{-1/2} \mathbf{x}_i$ , then  $\mathbf{y}_1, \dots, \mathbf{y}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ .

$$\begin{aligned} g(\boldsymbol{\theta}; \mathbf{v}) &= \frac{\frac{1}{n} \sum_{i=1}^n \left( \mathbf{r}^\top(\boldsymbol{\theta}) \mathbf{x}_i \right)^2 \left( (\nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{x}_i \right)^2}{\frac{1}{n} \sum_{i=1}^n \left( \mathbf{r}^\top(\boldsymbol{\theta}) \mathbf{x}_i \right)^2 \cdot \frac{1}{n} \sum_{i=1}^n \left( (\nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{x}_i \right)^2} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n \left( (\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}))^\top \mathbf{y}_i \right)^2 \left( (\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{y}_i \right)^2}{\frac{1}{n} \sum_{i=1}^n \left( (\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}))^\top \mathbf{y}_i \right)^2 \cdot \frac{1}{n} \sum_{i=1}^n \left( (\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{y}_i \right)^2}, \end{aligned}$$

588 Case(i). If  $\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}) = \mathbf{0}$  or  $\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v} = \mathbf{0}$ , we have  $g(\boldsymbol{\theta}; \mathbf{v}) = \frac{0}{0} = 1$ , this theorem holds.

589 Case (ii). If  $\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}) \neq \mathbf{0}$  and  $\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v} \neq \mathbf{0}$ , we define the following normalized vectors:

$$\tilde{\mathbf{r}}(\boldsymbol{\theta}) := \frac{\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta})}{\left\| \mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}) \right\|} \in \mathbb{S}^{d-1} \quad \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v}) := \frac{\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v}}{\left\| \mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v} \right\|} \in \mathbb{S}^{d-1}.$$

590 From the homogeneity of  $g(\boldsymbol{\theta}; \mathbf{v})$ , we have:

$$g(\boldsymbol{\theta}; \mathbf{v}) = \frac{\frac{1}{n} \sum_{i=1}^n \left( \tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \mathbf{y}_i \right)^2 \left( \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \mathbf{y}_i \right)^2}{\frac{1}{n} \sum_{i=1}^n \left( \tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \mathbf{y}_i \right)^2 \cdot \frac{1}{n} \sum_{i=1}^n \left( \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \mathbf{y}_i \right)^2}.$$

591 One the one hand, with the help of Lemma H.4, there exists  $C_1 > 0$  such that if we choose  
592  $n \gtrsim d + \log(1/\delta)$ , then w.p. at least  $1 - \delta/2$ , we have:

$$\inf_{\mathbf{w}, \mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 (\mathbf{u}^\top \mathbf{y}_i)^2 \geq C_1.$$

593 On the other hand, with the help of Lemma H.3, if we choose  $\epsilon = 1/2$  and  $n \gtrsim d + \log(1/\delta)$ , then  
 594 w.p. at least  $1 - \delta/2$ , we have:

$$\sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 \geq 1 + \frac{1}{2} = \frac{3}{2},$$

595 Combining these two bounds, we obtain that: if we choose  $\epsilon = 1/2$  and  $n \gtrsim d + \log(1/\delta)$ , then  
 596 w.p. at least  $1 - \delta$ , we have:

$$\begin{aligned} & \inf_{\mathbf{w}, \mathbf{u} \in \mathbb{S}^{d-1}} \frac{\frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 (\mathbf{u}^\top \mathbf{y}_i)^2}{\frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 \cdot \frac{1}{n} \sum_{i=1}^n (\mathbf{u}^\top \mathbf{y}_i)^2} \\ & \geq \frac{\inf_{\mathbf{w}, \mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 (\mathbf{u}^\top \mathbf{y}_i)^2}{\left( \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 \right)^2} \geq \frac{4C_1}{9}, \end{aligned}$$

597 which implies that

$$\inf_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p} g(\boldsymbol{\theta}; \mathbf{v}) \geq \min \left\{ 1, \inf_{\mathbf{w}, \mathbf{u} \in \mathbb{S}^{d-1}} \frac{\frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 (\mathbf{u}^\top \mathbf{y}_i)^2}{\frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 \cdot \frac{1}{n} \sum_{i=1}^n (\mathbf{u}^\top \mathbf{y}_i)^2} \right\} \geq \min \left\{ 1, \frac{4C_1}{9} \right\}.$$

598 □

## 599 H.2 Proof of Theorem 4.3

600 We first need a few lemmas.

601 **Lemma H.5.** Let  $\mathbf{y}_1, \dots, \mathbf{y}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . If  $n \gtrsim d^2 + \log^2(1/\delta)$ , then w.p. at least  $1 - \delta$ , we  
 602 have

$$\sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v})^4 \leq 8.$$

603 *Proof.* For  $\mathbb{S}^{d-1}$ , its covering number has the bound:

$$\left( \frac{1}{\rho} \right)^d \leq \mathcal{N}(\mathbb{S}^{d-1}, \rho) \leq \left( \frac{2}{\rho} + 1 \right)^d,$$

604 so there exist a  $\rho$ -net on  $\mathbb{S}^{d-1}$ :  $\mathcal{V} \subset \mathbb{S}^{d-1}$ , s.t.  $|\mathcal{V}| \leq \left( \frac{2}{\rho} + 1 \right)^d$ .

605 Step I. Bounding the term on the  $\rho$ -net.

606 For a fixed  $\mathbf{v} \in \mathcal{V}$ , due to  $\mathbf{y}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , we can verify  $(\mathbf{y}_i^\top \mathbf{v})^4$  is sub-Weibull random variable:

$$\mathbb{E} \exp \left( \left( (\mathbf{y}_i^\top \mathbf{v})^4 \right)^{1/2} \right) = \mathbb{E} \exp \left( (\mathbf{y}_i^\top \mathbf{v})^2 \right) \lesssim 1,$$

607 which means that there exist an absolute constant  $C_1 \geq 1$  s.t.  $\|(\mathbf{y}_i^\top \mathbf{v})^4\|_{\psi_{1/2}} \leq C_1$ .

608 By the concentration inequality for Sub-Weibull distribution with  $\beta = 1/2$  (Lemma J.5) and  
 609  $\mathbb{E}[(\mathbf{y}^\top \mathbf{v})^4] = 3$ , there exists an absolute constant  $C_2 \geq 1$  s.t.

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] - 3 \right| > \phi(n; \delta) \right) \leq 2\delta,$$

610 where  $\phi(n; \delta) = C_2 \left( \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log^2(1/\delta)}{n} \right)$ . Applying a union bound over  $\mathbf{v} \in \mathcal{V}$ , we have:

$$\mathbb{P} \left( \exists \mathbf{v} \in \mathcal{V} \text{ s.t. } \left| \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] - 3 \right| > \phi(n; \delta) \right)$$

$$\begin{aligned} &\leq \mathbb{P} \left( \bigcup_{\mathbf{v} \in \mathcal{V}} \left\{ \left| \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] - 3 \right| > \phi(n; \delta) \right\} \right) \leq \sum_{\mathbf{v} \in \mathcal{V}} \mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] - 3 \right| > \phi(n; \delta) \right) \\ &\leq 2|\mathcal{V}| \exp \left( -\frac{n}{C_2^2} \right) = 2 \left( \frac{2}{\rho} + 1 \right)^d \delta. \end{aligned}$$

611 So *w.p.* at least  $1 - 2 \left( \frac{2}{\rho} + 1 \right)^d \delta$ , we have:

$$\max_{\mathbf{v} \in \mathcal{V}} \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] \leq 3 + \phi(n; \delta).$$

612 Step II. Estimate the error of the  $\rho$ -net approximation.

613 For simplicity, we denote

$$P := \max_{\mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4], \quad Q := \max_{\mathbf{v} \in \mathcal{V}} \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4].$$

614 Let  $\mathbf{v} \in \mathbb{S}^{d-1}$  such that  $\frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] = P$ , then there exist  $\mathbf{v}_0 \in \mathcal{V}$ , s.t.  $\|\mathbf{v} - \mathbf{v}_0\| \leq \rho$ .

615 On the one hand,

$$\begin{aligned} &\left| \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v})^4 - \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v}_0)^4 \right| = \left| \frac{1}{n} \sum_{i=1}^n \left( (\mathbf{y}_i^\top \mathbf{v})^4 - (\mathbf{y}_i^\top \mathbf{v}_0)^4 \right) \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top (\mathbf{v} - \mathbf{v}_0)) (\mathbf{y}_i^\top (\mathbf{v} + \mathbf{v}_0)) \left( (\mathbf{y}_i^\top \mathbf{v})^2 + (\mathbf{y}_i^\top \mathbf{v}_0)^2 \right) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top (\mathbf{v} - \mathbf{v}_0)) (\mathbf{y}_i^\top (\mathbf{v} + \mathbf{v}_0)) (\mathbf{y}_i^\top \mathbf{v})^2 \right| + \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top (\mathbf{v} - \mathbf{v}_0)) (\mathbf{y}_i^\top (\mathbf{v} + \mathbf{v}_0)) (\mathbf{y}_i^\top \mathbf{v}_0)^2 \right| \\ &\leq \sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top (\mathbf{v} - \mathbf{v}_0))^2 (\mathbf{y}_i^\top (\mathbf{v} + \mathbf{v}_0))^2} \left( \sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v})^4} + \sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v}_0)^4} \right) \\ &\leq \sqrt[4]{\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top (\mathbf{v} - \mathbf{v}_0))^4} \sqrt[4]{\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top (\mathbf{v} + \mathbf{v}_0))^4} \left( \sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v})^4} + \sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v}_0)^4} \right) \\ &\leq \|\mathbf{v} - \mathbf{v}_0\| P^{1/4} \|\mathbf{v} + \mathbf{v}_0\| P^{1/4} (\sqrt{P} + \sqrt{Q}) \leq 2\rho \sqrt{P} (\sqrt{P} + \sqrt{Q}) \end{aligned}$$

616 On the other hand,

$$\left| \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v})^4 - \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v}_0)^4 \right| \geq P - \sum_{i=1}^n (\mathbf{y}_i^\top \mathbf{v}_0)^4 \geq P - Q.$$

617 Hence, we obtain

$$P - Q \leq 2\rho \sqrt{P} (\sqrt{P} + \sqrt{Q}),$$

618 which means that

$$P \leq \left( \frac{1}{1 - 2\rho} \right)^2 Q.$$

619 Step III. The bound for any  $\mathbf{v} \in \mathbb{S}^{d-1}$ .

620 Select  $\rho = \frac{1}{2} \left( 1 - \frac{1}{\sqrt{2}} \right)$  and denote  $\delta' = 2 \left( \frac{2}{\rho} + 1 \right)^d \delta$ . And we choose  $n \gtrsim d^2 + \log^2(1/\delta')$ , which  
621 ensures  $\phi(n; \delta) \leq 1$ .

622 Then combining the results in Step I and Step II, we know that: w.p. at least  $1 - \delta'$ , we have:

$$\max_{\mathbf{v} \in \mathcal{V}} \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] \leq 3 + 1 = 4; \quad \max_{\mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] \leq 2 \max_{\mathbf{v} \in \mathcal{V}} \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4],$$

623 which means

$$\max_{\mathbf{v} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n [(\mathbf{y}_i^\top \mathbf{v})^4] \leq 2 \cdot 4 = 8.$$

624

□

625 **Lemma H.6.** Let  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . For any  $\epsilon, \delta \in (0, 1)$ , if we choose

$$n \gtrsim \max \left\{ (d^2 \log^2(1/\epsilon) + \log^2(1/\delta)) / \epsilon, (d \log(1/\epsilon) + \log(1/\delta)) / \epsilon^2 \right\},$$

626 then w.p. at least  $1 - \delta$ , we have:

$$\sup_{\mathbf{w}, \mathbf{v} \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_1)^2 (\mathbf{v}^\top \mathbf{x}_1)^2] \right| \leq \epsilon.$$

627 *Proof.* For  $\mathbb{S}^{d-1}$ , its covering number has the bound:

$$\left( \frac{1}{\rho} \right)^d \leq \mathcal{N}(\mathbb{S}^{d-1}, \rho) \leq \left( \frac{2}{\rho} + 1 \right)^d,$$

628 so there exist two  $\rho$ -nets on  $\mathbb{S}^{d-1}$ :  $\mathcal{W} \subset \mathbb{S}^{d-1}$  and  $\mathcal{V} \subset \mathbb{S}^{d-1}$ , s.t.

$$|\mathcal{W}| \leq \left( \frac{2}{\rho} + 1 \right)^d, \quad |\mathcal{V}| \leq \left( \frac{2}{\rho} + 1 \right)^d.$$

629 Step I. Bounding the term on the  $\rho$ -net.

630 In this step, will estimate the term  $\left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E}[(\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2] \right|$  for any  $\mathbf{w} \in$   
631  $\mathcal{W}$  and  $\mathbf{v} \in \mathcal{V}$ .

632 For fixed  $\mathbf{w} \in \mathcal{W}$  and  $\mathbf{v} \in \mathcal{V}$ , we denote  $X_i^{\mathbf{w}, \mathbf{v}} := (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2$ . We can verify  $X_i$  is a  
633 sub-Weibull random variable with  $\beta = 1/2$  (Definition J.4):

$$\begin{aligned} & \mathbb{E} \left[ \exp \left( |(\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2|^{1/2} \right) \right] = \mathbb{E} \left[ \exp \left( |\mathbf{w}^\top \mathbf{x}_i| |\mathbf{v}^\top \mathbf{x}_i| \right) \right] \\ & \leq \mathbb{E} \left[ \exp \left( \frac{(\mathbf{w}^\top \mathbf{x}_i)^2 + (\mathbf{v}^\top \mathbf{x}_i)^2}{2} \right) \right] = \mathbb{E} \left[ \exp \left( \frac{(\mathbf{w}^\top \mathbf{x}_i)^2}{2} \right) \exp \left( \frac{(\mathbf{v}^\top \mathbf{x}_i)^2}{2} \right) \right] \\ & \stackrel{\text{Lemma J.6}}{\leq} \sqrt{\mathbb{E} \left[ \exp \left( (\mathbf{w}^\top \mathbf{x}_i)^2 \right) \right]} \cdot \sqrt{\mathbb{E} \left[ \exp \left( (\mathbf{v}^\top \mathbf{x}_i)^2 \right) \right]} \|\mathbf{v}^\top \mathbf{x}_i\|_{\psi_1}^2 \stackrel{\leq C_3}{\lesssim} 1, \end{aligned}$$

634 which means that there exists an absolute constant  $C_4 \geq 1$ , s.t.  $\|X_i^{\mathbf{w}, \mathbf{v}}\|_{\psi_{1/2}} \leq C_4$ . By the  
635 concentration inequality for Sub-Weibull distribution with  $\beta = 1/2$  (Lemma J.5), there exists an  
636 absolute constant  $C_5 \geq 1$ , s.t.

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n X_i^{\mathbf{w}, \mathbf{v}} - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i^{\mathbf{w}, \mathbf{v}}] \right| > \psi(n; \delta) \right) \leq \delta.$$

637 where  $\psi(n; \delta) = C_5 \left( \sqrt{\frac{\log(1/\delta)}{n}} + \frac{(\log(1/\delta))^2}{n} \right)$ .



638 Applying an union bound over  $\mathbf{w} \in \mathcal{W}$  and  $\mathbf{v} \in \mathcal{V}$ , we have:

$$\begin{aligned}
& \mathbb{P} \left( \exists \mathbf{w} \in \mathcal{W}, \mathbf{v} \in \mathcal{V}, \text{s.t.} \left| \frac{1}{n} \sum_{i=1}^n X_i^{\mathbf{w}, \mathbf{v}} - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i^{\mathbf{w}, \mathbf{v}}] \right| > \psi(n; \delta) \right) \\
& \leq \mathbb{P} \left( \bigcup_{(\mathbf{w}, \mathbf{v}) \in \mathcal{W} \times \mathcal{V}} \left\{ \exists \mathbf{w} \in \mathcal{W}, \mathbf{v} \in \mathcal{V}, \text{s.t.} \left| \frac{1}{n} \sum_{i=1}^n X_i^{\mathbf{w}, \mathbf{v}} - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i^{\mathbf{w}, \mathbf{v}}] \right| > \psi(n; \delta) \right\} \right) \\
& \leq \sum_{(\mathbf{w}, \mathbf{v}) \in \mathcal{W} \times \mathcal{V}} \mathbb{P} \left( \exists \mathbf{w} \in \mathcal{W}, \mathbf{v} \in \mathcal{V}, \text{s.t.} \left| \frac{1}{n} \sum_{i=1}^n X_i^{\mathbf{w}, \mathbf{v}} - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i^{\mathbf{w}, \mathbf{v}}] \right| > \psi(n; \delta) \right) \\
& \leq 2|\mathcal{W}||\mathcal{V}|\delta \leq 2 \left( \frac{2}{\rho} + 1 \right)^{2d} \delta.
\end{aligned}$$

639 So *w.p.* at least  $1 - 2 \left( \frac{2}{\rho} + 1 \right)^{2d} \delta$ , we have:

$$\sup_{\mathbf{w} \in \mathcal{W}, \mathbf{v} \in \mathcal{V}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \leq \psi(n; \delta).$$

640 Step II. Estimate the population error of the  $\rho$ -net approximation.

641 Let  $\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1}$ , s.t.  $\|\mathbf{w} - \mathbf{w}_0\| \leq \rho$  and  $\|\mathbf{v} - \mathbf{v}_0\| \leq \rho$ . For the population error, we have

$$\begin{aligned}
& \left| \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] - \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] \right| \\
& = \left| \mathbb{E} \left[ \left( (\mathbf{w}^\top \mathbf{x})^2 - (\mathbf{w}_0^\top \mathbf{x})^2 \right) (\mathbf{v}^\top \mathbf{x})^2 \right] + \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 \left( (\mathbf{v}^\top \mathbf{x})^2 - (\mathbf{v}_0^\top \mathbf{x})^2 \right) \right] \right| \\
& \leq \left| \mathbb{E} \left[ \left( (\mathbf{w}^\top \mathbf{x})^2 - (\mathbf{w}_0^\top \mathbf{x})^2 \right) (\mathbf{v}^\top \mathbf{x})^2 \right] \right| + \left| \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 \left( (\mathbf{v}^\top \mathbf{x})^2 - (\mathbf{v}_0^\top \mathbf{x})^2 \right) \right] \right|
\end{aligned}$$

642 We first bound  $\left| \mathbb{E} \left[ \left( (\mathbf{w}^\top \mathbf{x})^2 - (\mathbf{w}_0^\top \mathbf{x})^2 \right) (\mathbf{v}^\top \mathbf{x})^2 \right] \right|$ :

$$\begin{aligned}
& \left| \mathbb{E} \left[ \left( (\mathbf{w}^\top \mathbf{x})^2 - (\mathbf{w}_0^\top \mathbf{x})^2 \right) (\mathbf{v}^\top \mathbf{x})^2 \right] \right| = \left| \mathbb{E} \left[ \left( (\mathbf{w} - \mathbf{w}_0)^\top \mathbf{x} \mathbf{x}^\top (\mathbf{w} + \mathbf{w}_0) \right) (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \\
& \leq \left( \mathbb{E} \left[ \left( (\mathbf{w} - \mathbf{w}_0)^\top \mathbf{x} \mathbf{x}^\top (\mathbf{w} + \mathbf{w}_0) \right)^2 \right] \right)^{1/2} \left( \mathbb{E} \left[ (\mathbf{v}^\top \mathbf{x})^4 \right] \right)^{1/2} \\
& \leq \left( \mathbb{E} \left[ \left( (\mathbf{w} - \mathbf{w}_0)^\top \mathbf{x} \right)^4 \right] \right)^{1/4} \left( \mathbb{E} \left[ \left( (\mathbf{w} + \mathbf{w}_0)^\top \mathbf{x} \right)^4 \right] \right)^{1/4} \left( \mathbb{E} \left[ (\mathbf{v}^\top \mathbf{x})^4 \right] \right)^{1/2} \\
& \leq 3 \|\mathbf{w} - \mathbf{w}_0\| \|\mathbf{w} + \mathbf{w}_0\| \|\mathbf{v}\|^2 \leq 6\rho.
\end{aligned}$$

643 Repeating the proof above, we also have:

$$\left| \mathbb{E} \left[ \left( (\mathbf{w}^\top \mathbf{x})^2 - (\mathbf{w}_0^\top \mathbf{x})^2 \right) (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \leq 6\rho.$$

644 Combining these two inequalities, we have:

$$\left| \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] - \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] \right| \leq 6\rho + 6\rho = 12\rho.$$

645 Due to the arbitrariness of  $\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0$ , we obtain

$$\sup_{\substack{\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1} \\ \|\mathbf{w} - \mathbf{w}_0\| \leq \rho, \|\mathbf{v} - \mathbf{v}_0\| \leq \rho}} \left| \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] - \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] \right| \leq 12\rho.$$

646 Step III. Estimate the empirical error of the  $\rho$ -net approximation.

647 Let  $\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1}$ , s.t.  $\|\mathbf{w} - \mathbf{w}_0\| \leq \rho$  and  $\|\mathbf{v} - \mathbf{v}_0\| \leq \rho$ . For the empirical error, we have

$$\left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right|$$

$$\begin{aligned}
&= \left| \frac{1}{n} \sum_{i=1}^n \left[ \left( (\mathbf{w}^\top \mathbf{x}_i)^2 - (\mathbf{w}_0^\top \mathbf{x}_i)^2 \right) (\mathbf{v}^\top \mathbf{x}_i)^2 \right] + \frac{1}{n} \sum_{i=1}^n \left[ (\mathbf{w}_0^\top \mathbf{x}_i)^2 \left( (\mathbf{v}^\top \mathbf{x}_i)^2 - (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right) \right] \right| \\
&\leq \left| \frac{1}{n} \sum_{i=1}^n \left[ \left( (\mathbf{w}^\top \mathbf{x}_i)^2 - (\mathbf{w}_0^\top \mathbf{x}_i)^2 \right) (\mathbf{v}^\top \mathbf{x}_i)^2 \right] \right| + \left| \frac{1}{n} \sum_{i=1}^n \left[ (\mathbf{w}_0^\top \mathbf{x}_i)^2 \left( (\mathbf{v}^\top \mathbf{x}_i)^2 - (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right) \right] \right|
\end{aligned}$$

648 We first bound  $\left| \frac{1}{n} \sum_{i=1}^n \left[ \left( (\mathbf{w}^\top \mathbf{x}_i)^2 - (\mathbf{w}_0^\top \mathbf{x}_i)^2 \right) (\mathbf{v}^\top \mathbf{x}_i)^2 \right] \right|$ :

$$\begin{aligned}
\left| \frac{1}{n} \sum_{i=1}^n \left[ \left( (\mathbf{w}^\top \mathbf{x}_i)^2 - (\mathbf{w}_0^\top \mathbf{x}_i)^2 \right) (\mathbf{v}^\top \mathbf{x}_i)^2 \right] \right| &= \left| \frac{1}{n} \sum_{i=1}^n \left[ \left( (\mathbf{w} - \mathbf{w}_0)^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{w} + \mathbf{w}_0) \right) (\mathbf{v}^\top \mathbf{x}_i)^2 \right] \right| \\
&\leq 2\rho \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{u})^4.
\end{aligned}$$

649 Repeating the proof above, we also have  $\left| \frac{1}{n} \sum_{i=1}^n \left[ (\mathbf{w}_0^\top \mathbf{x}_i)^2 \left( (\mathbf{v}^\top \mathbf{x}_i)^2 - (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right) \right] \right| \leq$

650  $2\rho \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{u})^4$ . Combining these two bounds, we have:

$$\left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right| \leq 4\rho \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{u})^4.$$

651 Using Lemma H.5, if  $n \gtrsim d^2 + \log^2(1/\delta')$ , then w.p. at least  $1 - \delta'/2$ , we have

652  $\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{u})^4 \leq 8$ .

653 Hence, w.p. at least  $1 - \delta'/2$ , we have

$$\left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right| \leq 32\rho.$$

654 Due to the arbitrariness of  $\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0$ , we obtain

$$\sup_{\substack{\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1} \\ \|\mathbf{w} - \mathbf{w}_0\| \leq \rho, \|\mathbf{v} - \mathbf{v}_0\| \leq \rho}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right| \leq 32\rho.$$

655 Step IV. The bound for any  $\mathbf{w}, \mathbf{v} \in \mathbb{S}^{d-1}$ .

656 Combining the results in Step I, II, and II, we know that w.p. at least  $1 - \frac{\delta'}{2} - (\frac{2}{\rho} + 1)^d$ , we have

$$\begin{aligned}
\sup_{\mathbf{w} \in \mathcal{W}, \mathbf{v} \in \mathcal{V}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \right| &\leq \psi(n; \delta), \\
\sup_{\substack{\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1} \\ \|\mathbf{w} - \mathbf{w}_0\| \leq \rho, \|\mathbf{v} - \mathbf{v}_0\| \leq \rho}} \left| \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] - \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] \right| &\leq 12\rho, \\
\sup_{\substack{\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1} \\ \|\mathbf{w} - \mathbf{w}_0\| \leq \rho, \|\mathbf{v} - \mathbf{v}_0\| \leq \rho}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right| &\leq 32\rho.
\end{aligned}$$

657 Then for any  $\mathbf{w}, \mathbf{v} \in \mathbb{S}^{d-1}$ , there exists  $\mathbf{w}_0 \in \mathcal{W}, \mathbf{v}_0 \in \mathcal{V}$  s.t.  $\|\mathbf{w} - \mathbf{w}_0\| \leq \rho$  and  $\|\mathbf{v} - \mathbf{v}_0\| \leq \rho$ , so

$$\begin{aligned}
&\left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \\
&= \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 + \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right|
\end{aligned}$$

$$\begin{aligned}
& - \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] + \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \Big| \\
\leq & \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right| \\
& + \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 - \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] \right| + \left| \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \\
\leq & \sup_{\substack{\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1} \\ \|\mathbf{w} - \mathbf{w}_0\| \leq \rho, \|\mathbf{v} - \mathbf{v}_0\| \leq \rho}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_0^\top \mathbf{x}_i)^2 (\mathbf{v}_0^\top \mathbf{x}_i)^2 \right| \\
& + \sup_{\mathbf{w} \in \mathcal{W}, \mathbf{v} \in \mathcal{V}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \\
& + \sup_{\substack{\mathbf{w}, \mathbf{v}, \mathbf{w}_0, \mathbf{v}_0 \in \mathbb{S}^{d-1} \\ \|\mathbf{w} - \mathbf{w}_0\| \leq \rho, \|\mathbf{v} - \mathbf{v}_0\| \leq \rho}} \left| \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] - \mathbb{E} \left[ (\mathbf{w}_0^\top \mathbf{x})^2 (\mathbf{v}_0^\top \mathbf{x})^2 \right] \right| \\
\leq & 32\rho + \psi(n; \delta) + 12\rho = 44\rho + \psi(n; \delta).
\end{aligned}$$

658 Due to the arbitrariness of  $\mathbf{w}, \mathbf{v}$ , we have

$$\sup_{\mathbf{w}, \mathbf{v} \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \leq 44\rho + \psi(n; \delta)$$

659 Select  $\rho = \frac{\epsilon}{66}$  and  $\delta'/2 = 2(1 + \frac{2}{\rho})^{2d}\delta$ . And we choose

$$n \gtrsim \max \left\{ (d^2 \log^2(1/\epsilon) + \log^2(1/\delta)) / \epsilon, (d \log(1/\epsilon) + \log(1/\delta)) / \epsilon^2 \right\},$$

660 which satisfies  $\psi(n; \delta) \leq \epsilon/3$ .

661 Then *w.p.* at least  $1 - \delta'/2 - \delta'/2 = 1 - \delta'$ , we have

$$\sup_{\mathbf{w}, \mathbf{v} \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 (\mathbf{v}^\top \mathbf{x}_i)^2 - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{x})^2 (\mathbf{v}^\top \mathbf{x})^2 \right] \right| \leq \frac{44}{66}\epsilon + \frac{1}{3}\epsilon = \epsilon.$$

662

□

663 With the preparation of Lemma H.1, H.3, and H.6, now we give the proof of Theorem 4.3.

664 **Proof of Theorem 4.3.** Let  $\mathbf{y}_i = \mathbf{S}^{-1/2} \mathbf{x}_i$ , then  $\mathbf{y}_1, \dots, \mathbf{y}_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I_d)$ .

$$\begin{aligned}
g(\boldsymbol{\theta}; \mathbf{v}) &= \frac{\frac{1}{n} \sum_{i=1}^n \left( \mathbf{r}^\top(\boldsymbol{\theta}) \mathbf{x}_i \right)^2 \left( (\nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{x}_i \right)^2}{\frac{1}{n} \sum_{i=1}^n \left( \mathbf{r}^\top(\boldsymbol{\theta}) \mathbf{x}_i \right)^2 \cdot \frac{1}{n} \sum_{i=1}^n \left( (\nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{x}_i \right)^2} \\
&= \frac{\frac{1}{n} \sum_{i=1}^n \left( (\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}))^\top \mathbf{y}_i \right)^2 \left( (\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{y}_i \right)^2}{\frac{1}{n} \sum_{i=1}^n \left( (\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}))^\top \mathbf{y}_i \right)^2 \cdot \frac{1}{n} \sum_{i=1}^n \left( (\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{y}_i \right)^2},
\end{aligned}$$

665 Case (i). If  $\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}) = \mathbf{0}$  or  $\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v} = \mathbf{0}$ , we have  $g(\boldsymbol{\theta}; \mathbf{v}) = \frac{0}{0} = 1$ , this theorem holds.

666 Case (ii). If  $\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}) \neq \mathbf{0}$  and  $\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v} \neq \mathbf{0}$ , we define the following normalized vectors:

$$\tilde{\mathbf{r}}(\boldsymbol{\theta}) := \frac{\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta})}{\left\| \mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}) \right\|} \in \mathbb{S}^{d-1} \quad \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v}) := \frac{\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v}}{\left\| \mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v} \right\|} \in \mathbb{S}^{d-1}.$$

667 From the homogeneity of  $g(\boldsymbol{\theta}; \mathbf{v})$ , we have:

$$g(\boldsymbol{\theta}; \mathbf{v}) = \frac{\frac{1}{n} \sum_{i=1}^n \left( \tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \mathbf{y}_i \right)^2 \left( \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \mathbf{y}_i \right)^2}{\frac{1}{n} \sum_{i=1}^n \left( \tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \mathbf{y}_i \right)^2 \cdot \frac{1}{n} \sum_{i=1}^n \left( \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \mathbf{y}_i \right)^2}.$$

668 By Lemma H.3 and H.6, for any  $\epsilon, \delta \in (0, 1)$ , if we choose

$$n \gtrsim \max \left\{ \left( d^2 \log^2(1/\epsilon) + \log^2(1/\delta) \right) / \epsilon, \left( d \log(1/\epsilon) + \log(1/\delta) \right) / \epsilon^2 \right\},$$

669 then *w.p.* at least  $1 - \delta$ , the following inequalities hold:

$$\sup_{\mathbf{v} \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \mathbf{y}_i)^2 - 1 \right| \leq \epsilon,$$

$$\sup_{\mathbf{w}, \mathbf{v} \in \mathbb{S}^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{y}_i)^2 (\mathbf{v}^\top \mathbf{y}_i)^2 - \mathbb{E} \left[ (\mathbf{w}^\top \mathbf{y}_1)^2 (\mathbf{v}^\top \mathbf{y}_1)^2 \right] \right| \leq \epsilon;$$

670 These imply that for any  $\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p$ , we have:

$$\frac{\mathbb{E} \left[ (\tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \mathbf{y})^2 (\tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \mathbf{y})^2 \right] - \epsilon}{(1 + \epsilon)^2} \leq g(\boldsymbol{\theta}; \mathbf{v}) \leq \frac{\mathbb{E} \left[ (\tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \mathbf{y}_1)^2 (\tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \mathbf{y}_1)^2 \right] + \epsilon}{(1 - \epsilon)^2}. \quad (4)$$

671 First, we derive the upper bound for (4):

$$\begin{aligned} \text{RHS} &= \frac{\epsilon}{(1 - \epsilon)^2} + \frac{\mathbb{E} \left[ (\tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \mathbf{y})^2 (\tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \mathbf{y})^2 \right]}{(1 - \epsilon)^2 \left( \tilde{\mathbf{r}}(\boldsymbol{\theta})^\top \tilde{\mathbf{r}}(\boldsymbol{\theta}) \right) \left( \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v})^\top \tilde{\mathbf{w}}(\boldsymbol{\theta}; \mathbf{v}) \right)} \\ &\stackrel{\text{Homogeneity}}{=} \frac{\epsilon}{(1 - \epsilon)^2} + \frac{\mathbb{E} \left[ \left( (\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}))^\top \mathbf{y} \right)^2 \left( (\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v})^\top \mathbf{y} \right)^2 \right]}{(1 - \epsilon)^2 \left( (\mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}))^\top \mathbf{S}^{1/2} \mathbf{r}(\boldsymbol{\theta}) \right) \left( (\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v})^\top (\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v}) \right)} \\ &= \frac{\epsilon}{(1 - \epsilon)^2} + \frac{\mathbf{v}^\top \boldsymbol{\Sigma}(\boldsymbol{\theta}) \mathbf{v}}{2(1 - \epsilon)^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}} \stackrel{\text{Lemma H.1}}{=} \frac{\epsilon}{(1 - \epsilon)^2} + \frac{2\mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v} + (\nabla \mathcal{L}(\boldsymbol{\theta})^\top \mathbf{v})^2}{2(1 - \epsilon)^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}} \\ &= \frac{1 + \epsilon}{(1 - \epsilon)^2} + \frac{(\nabla \mathcal{L}(\boldsymbol{\theta})^\top \mathbf{v})^2}{2(1 - \epsilon)^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}} \stackrel{\text{Lemma H.2}}{\leq} \frac{1 + \epsilon}{(1 - \epsilon)^2} + \frac{1}{(1 - \epsilon)^2} = \frac{2 + \epsilon}{(1 - \epsilon)^2}. \end{aligned}$$

672 Moreover, if  $\langle \mathbf{v}, \mathcal{L}(\boldsymbol{\theta}) \rangle = 0$ , then the bound is

$$\text{RHS} \leq \frac{1 + \epsilon}{(1 - \epsilon)^2}.$$

673 In the similar way, we can derive the lower bound for (4):

$$\begin{aligned} \text{LHS} &= \frac{\mathbf{v}^\top \boldsymbol{\Sigma}(\boldsymbol{\theta}) \mathbf{v}}{2(1 + \epsilon)^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}} - \frac{\epsilon}{(1 + \epsilon)^2} \stackrel{\text{Lemma H.1}}{=} \frac{2\mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v} + (\nabla \mathcal{L}(\boldsymbol{\theta})^\top \mathbf{v})^2}{2(1 + \epsilon)^2 \mathcal{L}(\boldsymbol{\theta}) \mathbf{v}^\top G(\boldsymbol{\theta}) \mathbf{v}} - \frac{\epsilon}{(1 + \epsilon)^2} \\ &\geq \frac{1}{(1 + \epsilon)^2} - \frac{\epsilon}{(1 + \epsilon)^2} = \frac{1 - \epsilon}{(1 + \epsilon)^2}. \end{aligned}$$

674 So for any  $\mathbf{S}^{1/2} \mathbf{u}(\boldsymbol{\theta}) \neq \mathbf{0}$ ,  $\mathbf{S}^{1/2} \nabla F(\boldsymbol{\theta}) \mathbf{v} \neq \mathbf{0}$ , we have

$$\frac{1 - \epsilon}{(1 + \epsilon)^2} \leq g(\boldsymbol{\theta}; \mathbf{v}) \leq \frac{2 + \epsilon}{(1 - \epsilon)^2}.$$

675 Moreover, if  $\langle \mathbf{v}, \nabla \mathcal{L}(\boldsymbol{\theta}) \rangle = 0$ , then

$$\frac{1 - \epsilon}{(1 + \epsilon)^2} \leq g(\boldsymbol{\theta}; \mathbf{v}) \leq \frac{1 + \epsilon}{(1 - \epsilon)^2}.$$

676 Hence, we have proved this theorem: For any  $\epsilon, \delta > 0$ , if  $n \gtrsim$   
677  $\max \{ (d^2 \log^2(1/\epsilon) + \log^2(1/\delta)) / \epsilon, (d \log(1/\epsilon) + \log(1/\delta)) / \epsilon^2 \}$ , then *w.p.* at least  $1 - \delta$ ,  
678 the strong alignment holds uniformly:

$$\begin{aligned} \text{(i). } & \frac{1 - \epsilon}{(1 + \epsilon)^2} \leq \inf_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p} g(\boldsymbol{\theta}; \mathbf{v}) \leq \sup_{\boldsymbol{\theta}, \mathbf{v} \in \mathbb{R}^p} g(\boldsymbol{\theta}; \mathbf{v}) \leq \frac{2 + \epsilon}{(1 - \epsilon)^2}, \\ \text{(ii). } & \frac{1 - \epsilon}{(1 + \epsilon)^2} \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^p, \langle \mathbf{v}, \nabla \mathcal{L}(\boldsymbol{\theta}) \rangle = 0} g(\boldsymbol{\theta}; \mathbf{v}) \leq \sup_{\boldsymbol{\theta} \in \mathbb{R}^p, \langle \mathbf{v}, \nabla \mathcal{L}(\boldsymbol{\theta}) \rangle = 0} g(\boldsymbol{\theta}; \mathbf{v}) \leq \frac{1 + \epsilon}{(1 - \epsilon)^2}. \end{aligned}$$

679 □

## 680 I Proofs in Section 5: Escape directions

### 681 I.1 Proof of Theorem 5.2

682 Recall that  $\mathbf{w}(t) = \sum_{i=1}^d w_i(t) \mathbf{u}_i$  with  $w_i(t) = \mathbf{u}_i^\top \mathbf{w}(t)$ . Then,  $w_i(t+1) = (1 - \eta \lambda_i) w_i(t) +$   
683  $\eta \boldsymbol{\xi}(t)^\top \mathbf{u}_i$ . Taking the expectation of the square of both sides, we obtain

$$\mathbb{E}[w_i^2(t+1)] = (1 - \eta \lambda_i)^2 \mathbb{E}[w_i^2(t)] + \eta^2 \mathbb{E}[|\mathbf{u}_i^\top \boldsymbol{\xi}(t)|^2],$$

684 According to Assumption 5.1, there exists  $A_1, A_2 > 0$  such that for any  $i \in [d]$ ,

$$A_1 \lambda_i \mathcal{L}(\mathbf{w}_t) \leq \mathbb{E}[|\mathbf{u}_i^\top \boldsymbol{\xi}(t)|] \leq A_2 \lambda_i \mathcal{L}(\mathbf{w}_t).$$

685 Let  $X_t = \sum_{i=1}^k \lambda_i \mathbb{E}[w_i^2(t)]$ ,  $Y_t = \sum_{i=k+1}^d \lambda_i \mathbb{E}[w_i^2(t)]$  denote the components of loss energy along  
686 sharp and flat directions, respectively. And we denote  $D_k(t) := Y_t / X_t$ .

687 Plugging the fact that  $2\mathcal{L}(\mathbf{w}(t)) = X_t + Y_t$  into the two formulations above, we can obtain the  
688 following component dynamics:

$$\begin{aligned} X_{t+1} & \leq \alpha_k X_t + A_2 \eta^2 \left( \sum_{i=1}^k \lambda_i^2 \right) (X_t + Y_t), \\ X_{t+1} & \geq A_1 \eta^2 \left( \sum_{i=1}^k \lambda_i^2 \right) (X_t + Y_t), \\ Y_{t+1} & \geq A_1 \eta^2 \left( \sum_{i=k+1}^d \lambda_i^2 \right) (X_t + Y_t), \end{aligned} \tag{5}$$

689 where  $\alpha_k \leq \max_{i=1, \dots, k} |1 - \eta \lambda_i|^2$ . The terms  $\alpha_k X_t$  and  $\beta_k Y_t$  capture the impact of the gradient,  
690 while the remaining terms originate from the noise.

691 From (5), we have the following estimate about  $D_k(t+1)$ :

$$\begin{aligned} D_k(t+1) & = \frac{Y_{t+1}}{X_{t+1}} \geq \frac{A_1 \eta^2 \left( \sum_{i=k+1}^d \lambda_i^2 \right) (X_t + Y_t)}{\alpha_k X_t + A_2 \eta^2 \left( \sum_{i=1}^k \lambda_i^2 \right) (X_t + Y_t)} \\ & = \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{\alpha_k X_t}{A_2 \eta^2 \sum_{i=k+1}^d \lambda_i^2 (X_t + Y_t)}} \\ & \geq \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{\max_{1 \leq i \leq k} |1 - \eta \lambda_i|^2 X_t}{A_2 \eta^2 \sum_{i=1}^k \lambda_i^2 (X_t + Y_t)}}. \end{aligned} \tag{6}$$

692 We will prove this theorem for the learning rate  $\eta = \frac{\beta}{\|\mathcal{G}(\boldsymbol{\theta})\|_F}$ , where  $\beta \geq \frac{1.1}{\sqrt{A_1}}$ .

693 Case (I). Small learning rate  $\eta \in \left[ \frac{1.1}{\sqrt{A_1} \|\mathcal{G}(\boldsymbol{\theta})\|_F}, \frac{1}{\lambda_1} \right]$ .

694 In this step, we consider  $\eta = \frac{\beta}{\|G(\theta)\|_F}$  such that  $\beta \geq \frac{1.1}{\sqrt{A_1}}$  and  $\eta \leq \frac{1}{\lambda_1}$ . Then we have:

$$\frac{\max_{1 \leq i \leq k} |1 - \eta \lambda_i|^2}{A_2 \eta^2 \sum_{i=k+1}^d \lambda_i^2} \leq \frac{1}{A_2 \eta^2 \sum_{i=1}^k \lambda_i^2}.$$

695 Notice that (5) also ensures:

$$(X_{t+1} + Y_{t+1}) \geq A_1 \eta^2 \left( \sum_{i=1}^d \lambda_i^2 \right) (X_t + Y_t).$$

696 Combining this inequality with (5), we have the estimate:

$$\begin{aligned} \frac{X_{t+1}}{X_{t+1} + Y_{t+1}} &\leq \frac{\alpha_k X_t + A_2 \eta^2 \left( \sum_{i=1}^k \lambda_i^2 \right) (X_t + Y_t)}{X_{t+1} + Y_{t+1}} \\ &\leq \frac{\alpha_k X_t}{A_1 \eta^2 \left( \sum_{i=1}^d \lambda_i^2 \right) (X_t + Y_t)} + \frac{A_2 \left( \sum_{i=1}^k \lambda_i^2 \right)}{A_1 \left( \sum_{i=1}^d \lambda_i^2 \right)} \end{aligned}$$

697 For simplicity, we denote  $W_t := \frac{X_t}{X_t + Y_t}$ ,  $A := \frac{\alpha_k}{A_1 \eta^2 \left( \sum_{i=1}^d \lambda_i^2 \right)}$ , and  $B := \frac{A_2 \left( \sum_{i=1}^k \lambda_i^2 \right)}{A_1 \left( \sum_{i=1}^d \lambda_i^2 \right)}$ .

698 From  $\eta \leq 1/3$ , we have  $\alpha_k \leq 1$  and  $A \leq \frac{1}{A_1 \eta^2 \left( \sum_{i=1}^d \lambda_i^2 \right)} = \frac{1}{A_1 \beta^2} < 1$ . Moreover, it holds that

$$\begin{aligned} W_{t+1} &\leq A W_t + B \leq A(A W_{t-1} + B) + B = A^2 W_{t-1} + B(1 + A) \\ &\leq \dots \leq A^{t+1} W_0 + B(1 + A + \dots + A^t) = A^{t+1} W_0 + \frac{1 - A^{t+1}}{1 - A} B \end{aligned}$$

699 On the one hand, if we choose

$$t \geq \frac{\log \left( 1/W_0 A_2 \eta^2 \sum_{i=1}^k \lambda_i^2 \right)}{\log(A_1 \beta^2)},$$

700 then we have

$$A^t W_0 \leq \left( \frac{\alpha_k}{A_1 \eta^2 \left( \sum_{i=1}^d \lambda_i^2 \right)} \right)^t W_0 \leq \left( \frac{1}{A_1 \beta^2} \right)^t W_0 \leq A_2 \eta^2 \sum_{i=1}^k \lambda_i^2.$$

701 On the other hand, if we choose  $t \geq 1$ , then it holds that

$$\frac{1 - A^t}{1 - A} B \leq B = \frac{A_2 \left( \sum_{i=1}^k \lambda_i^2 \right)}{A_1 \left( \sum_{i=1}^d \lambda_i^2 \right)} \leq A_2 \eta^2 \sum_{i=1}^k \lambda_i^2.$$

702 Hence, if we choose

$$t \geq \max \left\{ 1, \frac{\log \left( 1/W_0 A_2 \eta^2 \sum_{i=1}^k \lambda_i^2 \right)}{\log(A_1 \beta^2)} \right\},$$

703 then we have

$$\frac{X_t}{X_t + Y_t} = W_t \leq A^t W_0 + \frac{1 - A^t}{1 - A} B \leq 2 A_2 \eta^2 \sum_{i=1}^k \lambda_i^2,$$

704 which implies that

$$\text{RHS of (6)} \geq \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{\max_{1 \leq i \leq k} |1 - \eta \lambda_i|^2}{A_2 \eta^2 \sum_{i=1}^k \lambda_i^2} \frac{X_t}{X_t + Y_t}}$$

$$\geq \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{1}{A_2 \eta^2 \sum_{i=1}^k \lambda_i^2}} \cdot 2A_2 \eta^2 \sum_{i=1}^k \lambda_i^2 = \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{3A_2 \sum_{i=1}^k \lambda_i^2}.$$

705 **Case (II).** Large learning rate  $\eta \geq 1/\lambda_1$ .

706 In this step, we consider  $\eta \geq \frac{1}{\lambda_1}$ . Then for any  $t \geq 0$ , we have:

$$\begin{aligned} \text{RHS of (6)} &= \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{\alpha_k}{\sum_{i=k+1}^d \lambda_i^2} \frac{X_t}{X_t + Y_t}} \geq \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{\max_{i \in [k]} |1 - \eta \lambda_i|^2}{A_2 \eta^2 \sum_{i=1}^k \lambda_i^2}} \\ &\geq \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{\max\{1, |1 - \eta \lambda_1|^2\}}{A_2 \eta^2 \sum_{i=1}^k \lambda_i^2}} \geq \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{A_2 \sum_{i=1}^k \lambda_i^2} \cdot \frac{1}{1 + \frac{1}{A_2}} = \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{(A_2 + 1) \sum_{i=1}^k \lambda_i^2}. \end{aligned}$$

707 Combining Case (I) and (II), we obtain this theorem: If we choose the learning rate  $\eta = \frac{\beta}{\|G(\theta)\|_F}$ ,

708 where  $\beta \geq \frac{1.1}{\sqrt{A_1}}$ , then for any

$$t \geq \max \left\{ 1, \frac{\log \left( 1/W_0 A_2 \eta^2 \sum_{i=1}^k \lambda_i^2 \right)}{\log(A_1 \beta^2)} \right\},$$

709 we have

$$D_k(t+1) \geq \frac{A_1 \sum_{i=k+1}^d \lambda_i^2}{\max\{3A_2, A_2 + 1\} \sum_{i=1}^k \lambda_i^2}.$$

710

□

## 711 I.2 Proof of Proposition 5.3

712 Recall that  $\mathbf{w}(t) = \sum_{i=1}^d w_i(t) \mathbf{u}_i$  with  $w_i(t) = \mathbf{u}_i^\top \mathbf{w}(t)$ . Then, for GD,  $w_i(t+1) = (1 - \eta \lambda_i) w_i(t)$ ,

713 which implies:

$$w_i(t) = (1 - \eta \lambda_i)^t w_i(0).$$

714 Therefore, for  $\eta = \beta/\lambda_1$  ( $\beta > 2$ ), it holds that

$$D_1(t) = \frac{\sum_{i=2}^d \lambda_i w_i^2(t)}{\lambda_1 w_1^2(t)} = \frac{\sum_{i=2}^d \lambda_i (1 - \eta \lambda_i)^{2t} w_i^2(0)}{\lambda_1 (1 - \eta \lambda_1)^{2t} w_1^2(0)}.$$

715

□

## 716 J Useful Inequalities

717 **Lemma J.1** (Bernstein's Inequality (Vershynin, 2018)). Suppose  $\{X_1, \dots, X_n\}$  are independent  
718 sub-Exponential random variables with  $\|X_i\|_{\psi_1} \leq K$ . Then there exists an absolute constant  $c > 0$   
719 such that for any  $t \geq 0$ , we have:

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] \right| > t \right) \leq 2 \exp \left( -cn \min \left\{ \frac{t}{K}, \frac{t^2}{K^2} \right\} \right).$$

720 **Lemma J.2** (Hanson-Wright's Inequality (Vershynin, 2018)). Let  $\mathbf{X} = (X_1, \dots, X_n) \in \mathbb{R}^n$  be a  
721 random vector with independent mean zero sub-Gaussian coordinates. Let  $\mathbf{A}$  be an  $n \times n$  matrix.  
722 Then, there exists an absolute constant  $c$  such that for every  $t \geq 0$ , we have

$$\mathbb{P} \left( \left| \mathbf{X}^\top \mathbf{A} \mathbf{X} - \mathbb{E}[\mathbf{X}^\top \mathbf{A} \mathbf{X}] \right| \geq t \right) \leq 2 \exp \left( -c \min \left\{ \frac{t^2}{K^4 \|\mathbf{A}\|_F^2}, \frac{t}{K^2 \|\mathbf{A}\|_2} \right\} \right),$$

723 where  $K = \max_i \|X_i\|_{\psi_2}$ .



724 **Lemma J.3** (Covariance Estimate for sub-Gaussian Distribution (Vershynin, 2018)). *Let*  
725  $\mathbf{x}, \mathbf{x}_1, \dots, \mathbf{x}_n$  *be i.i.d. random vectors in*  $\mathbb{R}^d$ . *More precisely, assume that there exists*  $K \geq 1$   
726 *s.t.  $\|\langle \mathbf{x}, \mathbf{v} \rangle\|_{\psi_2} \leq K \|\langle \mathbf{x}, \mathbf{v} \rangle\|_{L_2}$  for any  $\mathbf{v} \in \mathbb{S}^{d-1}$ , Then for any  $u \geq 0$ , w.p. at least  $1 - 2 \exp(-u)$*   
727 *one has*

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \right\| \leq CK^2 \left( \sqrt{\frac{d+u}{n}} + \frac{d+u}{n} \right) \|\mathbb{E}[\mathbf{x} \mathbf{x}^\top]\|,$$

728 *where*  $C$  *is an absolute positive constant.*

729 **Definition J.4** (Sub-Weibull Distribution). We define  $X$  as a sub-Weibull random variable if it has a  
730 bounded  $\psi_\beta$ -norm. The  $\psi_\beta$ -norm of  $X$  for any  $\beta > 0$  is defined as

$$\|X\|_{\psi_\beta} := \inf \left\{ C > 0 : \mathbb{E}[\exp(|X|^\beta / C^\beta)] \leq 2 \right\}.$$

731 Particularly, when  $\beta = 1$  or  $2$ , sub-Weibull random variables reduce to sub-Exponential or sub-  
732 Gaussian random variables, respectively.

733 **Lemma J.5** (Concentration Inequality for Sub-Weibull Distribution, Theorem 3.1 in (Hao et al.,  
734 2019)). *Suppose  $\{X_i\}_{i=1}^n$  are independent sub-Weibull random variables with  $\|X_i\|_{\psi_\beta} \leq K$ . Then*  
735 *there exists an absolute constant  $C_\beta$  only depending on  $\beta$  such that for any  $\delta \in (0, 1/e^2)$ , w.p. at*  
736 *least  $1 - \delta$ , we have*

$$\left| \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] \right| \leq C_\beta K \left( \left( \frac{\log(1/\delta)}{n} \right)^{1/2} + \frac{(\log(1/\delta))^{1/\beta}}{n} \right).$$

737 **Lemma J.6** (Cauchy-Schwarz Inequalities).

738 (1) *Let*  $S \in \mathbb{R}^{n \times n}$  *be a positive symmetric definite matrix. For any*  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , *we denote*  $\langle \mathbf{x}, \mathbf{y} \rangle_S :=$   
739  $\mathbf{x}^\top S \mathbf{y}$  *and*  $\|\mathbf{x}\|_S := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_S}$ , *then we have*  $|\langle \mathbf{x}, \mathbf{y} \rangle_S| \leq \|\mathbf{x}\|_S \|\mathbf{y}\|_S$ .

740 (2) *Given two random variables*  $X$  *and*  $Y$ , *it holds that*  $|\mathbb{E}[XY]| \leq \sqrt{\mathbb{E}[X^2]} \sqrt{\mathbb{E}[Y^2]}$ .