
Blessing of Depth in Linear Regression: Deeper Models Have Flatter Landscape Around the True Solution

Jianhao Ma

Department of Industrial & Operations Engineering
University of Michigan
Ann Arbor, MI 48109
jianhao@umich.edu

Salar Fattahi

Department of Industrial & Operations Engineering
University of Michigan
Ann Arbor, MI 48109
fattahi@umich.edu

Abstract

This work characterizes the effect of depth on the optimization landscape of linear regression, showing that, despite their nonconvexity, deeper models have more desirable optimization landscape. We consider a robust and over-parameterized setting, where a subset of measurements are grossly corrupted with noise, and the true linear model is captured via an N -layer diagonal linear neural network. On the negative side, we show that this problem *does not* have a benign landscape: given any $N \geq 1$, with constant probability, there exists a solution corresponding to the ground truth that is neither local nor global minimum. However, on the positive side, we prove that, for any N -layer model with $N \geq 2$, a simple sub-gradient method becomes oblivious to such “problematic” solutions; instead, it converges to a balanced solution that is not only close to the ground truth but also enjoys a flat local landscape, thereby eschewing the need for “early stopping”. Lastly, we empirically verify that the desirable optimization landscape of deeper models extends to other robust learning tasks, including deep matrix recovery and deep ReLU networks with ℓ_1 -loss.

1 Introduction

Supported by the empirical success of deep models in contemporary learning tasks, it is by now a conventional wisdom that “deeper models generalize better” [21, 31, 7]. Indeed, the flurry of recent attempts towards demystifying this phenomenon is a testament to the amount of research it has spawned: from simple linear regression to more complex and nonlinear models, it is shown that deeper models benefit from a range of desirable statistical properties, such as *depth separation* [33, 15, 34, 35], *implicit bias* [19, 2, 10], and *hierarchical learning* [1], to name a few.

Despite the great promise of deeper models—both theoretically and empirically—the effect of depth on their optimization landscape has remained elusive to this day. A recent body of work attempts to characterize the effect of depth on the loss function through the notion of *benign landscape*. Roughly speaking, an optimization problem has a benign landscape if it is devoid of spurious local minima, and its true solutions—i.e., solutions corresponding to the ground truth—coincide with global minima.

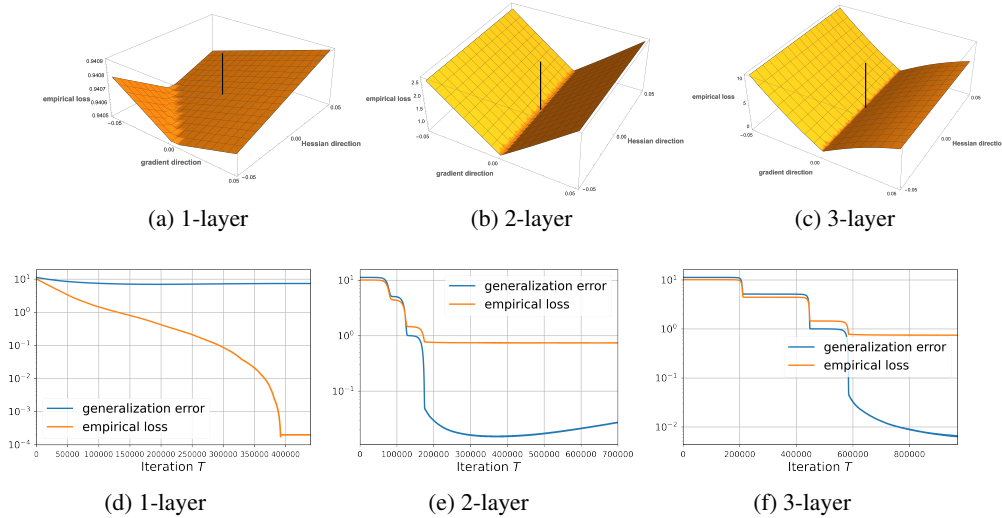


Figure 1: **First row.** Local landscape around the balanced true solution $\mathbf{w}^* = (R^{-\gamma}; \dots; R^{-\gamma})$ for 1-, 2-, and 3-layer models. x and y axis correspond to the points of the form $\mathbf{w}^* + \mathbf{d} + \mathbf{h}$, for different values of α and β , where \mathbf{d} and \mathbf{h} are respectively the most descent sub-gradient direction and the most negatively curved direction of the Hessian after smoothing. **Second row.** Generalization error and the empirical loss of the solutions found by SubGM for 1-, 2-, and 3-layer models.

It has been shown that 2-layer [5] and multi-layer [23] linear neural networks with nearly-noiseless data have benign landscape. However, the notion of benign landscape is significantly stronger than what is needed in practice. For instance, the existence of spurious local minima may not pose any issue if an optimization algorithm can avoid them efficiently. Another line of research focuses on characterizing the solution trajectory of different local-search algorithms, showing that they enjoy an *implicit bias* that steers them away from undesirable solutions [25, 38, 10, 19, 2, 18]. However, such guarantees only apply to specific trajectories of an algorithm, thereby falling short of any meaningful characterization of the optimization landscape around those trajectories.

1.1 Our Contributions

To shed light on the effect of depth on the optimization landscape of deep models, we consider a prototypical problem in machine learning, namely *robust linear regression*, where the goal is to recover a linear model from a limited number of grossly corrupted measurements. Given samples of the form $y_i = \langle \theta^*, x_i \rangle + \varepsilon_i$, we study the optimization landscape of ℓ_1 -loss under an N -layer model defined as $y = f_{\mathbf{w}}(x) := \langle w_1 \odot w_2 \odot \dots \odot w_N, x \rangle$. Our results are summarized as follows:

- We prove that, for any $N \geq 1$, there exists at least one true solution that is neither local nor global minimum of ℓ_1 -loss, provided that at least a fraction $p > 0$ of the measurements are corrupted with noise.
- Despite the ubiquity of such “hidden” true solutions, we show that, for any N -layer model with $N \geq 2$, a simple sub-gradient method (SubGM) with small initialization converges to a small neighborhood of a balanced true solution. The radius of this neighborhood shrinks with the depth of the model, resulting in more accurate solutions. Moreover, the balancedness of the solution implies that each layer of the model inherits a similar sparsity pattern to the ground truth.
- We prove that deeper models take longer to train, but once trained, the algorithm will stay close to the ground truth for a longer time. This implies that early stopping of the algorithm becomes less crucial for deeper models.
- Finally, we prove that depth flattens the optimization landscape around the solution obtained by SubGM. In particular, we show that, within an γ -neighborhood of the true solution, the steepest descent direction can reduce the loss by at most $\mathcal{O}(\gamma^N)$, which decreases *exponentially* with N .

Motivating Example. To showcase our results, we consider an instance of robust linear regression in the over-parameterized setting, where the dimension of $\theta^?$ is 500 and the number of available samples is 300. Moreover, we assume that 10% of the measurements are corrupted with large noise. The first row of Figure 1 shows the landscape around the balanced ground truth $w_1^? = \sqrt[N]{\theta^?}$, $w_2^? = \sqrt[N]{\theta^?}, \dots, w_N^? = \sqrt[N]{\theta^?}$.¹ In particular, x and y axis show the most descent sub-gradient direction and the most negatively curved direction of the Hessian after smoothing.² Evidently, there is a sharp transition in the landscape of N -layer models: for the 1-layer model, the true solution has strictly negative directions along both sub-gradient and negative curvature. However, these descent directions almost disappear in 2- and 3-layer models. The second row of Figure 1 shows the performance of SubGM on these models. It can be observed that a 1-layer model easily overfits to noise, leading to a vacuous generalization error. On the contrary, a 3-layer model can find a solution that generalizes progressively better than 1- and 2-layer models, demonstrating the algorithmic benefit of the depth.

Notations: For two vectors $x, y \in \mathbb{R}^d$, their inner product is defined as $\langle x, y \rangle = x^\top y$, and their Hadamard product is defined as $x \odot y = [x_1 y_1 \dots x_d y_d]^\top$. For simplicity of notation, we use $\prod_j w_j$ to denote the Hadamard product of $w_1, w_2, \dots, w_N \in \mathbb{R}^d$. For a vector x , $\|x\|$, $\|x\|_7$, and $\|x\|_0$ refer to 2-norm, ∞ -norm, and the number of nonzero elements, respectively. The symbols $a_t \cdot b_t$ and $a_t = \mathcal{O}(b_t)$ are used to denote $a_t \leq C b_t$, for a universal constant C . The notation $a_t = \Theta(b_t)$ is used to denote $a_t = \mathcal{O}(b_t)$ and $b_t = \mathcal{O}(a_t)$. The $\text{Sign}(\cdot)$ function is defined as $\text{Sign}(x) = x/|x|$ if $x \neq 0$, and $\text{Sign}(0) = [-1, 1]$. We denote $[n] := \{1, 2, \dots, n\}$. For a vector $x \in \mathbb{R}^d$, we define $x^a = [x_1^a x_2^a \dots x_d^a]^\top$, for any $a > 0$. In all of our probabilistic arguments, the randomness is only over the input data and noise.

2 Problem Formulation

We study the problem of robust and sparse linear regression, where the goal is to estimate a k -sparse vector $\theta^? \in \mathbb{R}^d$ ($k \ll d$) from a limited number of data points $\{(x_i, y_i)\}_{i=1}^m$, where $y_i = \langle \theta^?, x_i \rangle + \varepsilon_i$, x_i is i.i.d. standard Gaussian vector, and ε_i is noise. Moreover, for simplicity of our subsequent analysis, we assume that $\theta^?$ is a non-negative vector. We believe that this assumption can be relaxed without a significant change in our results.

Assumption 1 (Noise Model). *Given a corruption probability p , the noise vector $\varepsilon = [\varepsilon_1 \dots \varepsilon_m]^\top \in \mathbb{R}^m$ is generated as follows: first, a subset $\mathcal{S} \subset [m]$ with cardinality pm is chosen uniformly at random³. Then, for each entry $i \in \mathcal{S}$, the value of ε_i is drawn independently from a distribution P_o , and all the remaining entries are set to zero. Moreover, a random variable ζ under the distribution P_o satisfies $\mathbb{E}_{P_o}[\zeta] = 0$ and $\mathbb{P}(|\zeta| \geq t_0) \geq p_0$, for some strictly positive constants t_0 and p_0 .*

Our considered noise model does not impose any assumption on the magnitude of the noise or the specific form of its distribution, which makes it particularly suitable for modeling outliers. Note that the assumption $\mathbb{P}(|\varepsilon| \geq t_0) \geq p_0$ is very mild and satisfied for almost all common distributions. Roughly speaking, it implies that the noise takes a nonzero value with a nonzero probability.

To capture the input-output relationship, we consider a class of N -layer diagonal linear neural networks of the form $f_{\mathbf{w}}(x) = \langle w_1 \odot \dots \odot w_N, x \rangle$, where $\mathbf{w} := (w_1, \dots, w_N)$ collects the weights of the layers $w_1, \dots, w_N \in \mathbb{R}^d$. Due to the sparse-and-large nature of the noise, it is natural to minimize the so-called empirical risk with ℓ_1 -loss:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) := \frac{1}{m} \sum_{i=1}^m |f_{\mathbf{w}}(x_i) - y_i| = \frac{1}{m} \sum_{i=1}^m |\langle w_1 \odot \dots \odot w_N, x_i \rangle - y_i|. \quad (1)$$

Other variants of empirical risk minimization for linear regression have been studied in the literature. For instance, [10] study the solution trajectory of gradient flow on ℓ_2 -loss, showing that it converges to a solution with the smallest ℓ_1 -norm. Similar analysis has also appeared in more general deep linear neural networks [13, 14]. However, it is well-known that ℓ_2 -loss is highly sensitive to outliers, and ℓ_1 -loss is a better alternative to identify and reject large-and-sparse noise.

¹Later, we will show that a simple sub-gradient method converges to this balanced solution.

²To smoothen $|x|$, we replace it with $\sqrt{x^2 + \epsilon}$, for $\epsilon = 10^{-7}$.

³Here, for simplicity we assume pm is an integer.

A solution $\bar{\mathbf{w}}$ is called *global* if it corresponds to a global minimizer of $\mathcal{L}(\mathbf{w})$. Moreover, a *local solution* $\bar{\mathbf{w}}$ corresponds to the minimum of $\mathcal{L}(\mathbf{w})$ within an open ball centered at $\bar{\mathbf{w}}$. Finally, a *true solution* $\bar{\mathbf{w}}$ satisfies $\bar{w}_1 \odot \cdots \odot \bar{w}_N = \theta^\gamma$.

3 Main Results

3.1 Absence of Benign Landscape

We show that for any arbitrary corruption probability $0 < p < 1/2$ and any number of layers $N \geq 1$, there exists at least one true solution with a strictly negative descent direction, provided that the problem is *over-parameterized*, i.e., $m \gg d$.

Theorem 1 (unidentifiable true solutions). *Define $\mathcal{W} = \{\mathbf{w} : w_1 \odot \cdots \odot w_N = \theta^\gamma\}$ as the set of all true solutions of an N -layer model. For any $N \geq 1$ and $0 < p < 1/2$, the following statements hold:*

- (Over-parameterized regime) *If $m \leq 0.1d$, with probability of at least $1/16$, we have*

$$\inf_{\mathbf{w}^\gamma \in \mathcal{W}} \inf_{\mathbf{w} : k\mathbf{w}} \inf_{\mathbf{w}^\gamma \in k_1} \{\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}^\gamma)\} \leq -p_0 p \gamma, \quad (2)$$

for any $\gamma \in [t_0, 1]$.

- (Under-parameterized regime) *If $m \gg \frac{d}{(1-2p)^2}$, with probability of at least $1 - e^{-\Omega(d)}$, we have*

$$\inf_{\mathbf{w}^\gamma \in \mathcal{W}} \inf_{\mathbf{w}} \{\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}^\gamma)\} \geq 0. \quad (3)$$

The above proposition unravels a sharp transition in the landscape of robust linear regression with an N -layer model: when $m \gg d$, some of the true solutions are likely to be non-critical points, and hence, cannot be recovered via any first-order algorithm. As soon as $m \ll d$, all true solutions become global. This is in stark contrast with the recent results on the benign landscape of robust low-rank matrix recovery with ℓ_1 -loss, which show that, under the so-called *restricted isometry property* (RIP), all the true solutions are global and vice versa [26, 12, 16]. The vector version of RIP is known to hold with $m = \tilde{\Omega}(k)$ samples (see e.g. [3] for a simple proof). Theorem 1 shows that, unlike the low-rank matrix recovery, RIP is *not enough* to guarantee the equivalence between the true and global solutions in deep linear models.

The detailed proof of this theorem can be found in Appendix C.1. Here, we provide a proof sketch for $N = 1$ and $N = 2$ to elucidate the key ideas.

Proof sketch of Theorem 1. For 1-layer model, the set of true solutions \mathcal{W} reduces to a singleton $\{\theta^\gamma\}$. Upon choosing $\mathbf{w} = \theta^\gamma$, we prove the existence of a perturbation $\|\Delta\theta\| \leq \gamma$ that satisfies $\mathcal{L}(\theta^\gamma) - \mathcal{L}(\theta^\gamma + \Delta\theta) < 0$. The perturbed loss takes the following form

$$\mathcal{L}(\theta^\gamma + \Delta\theta) = \frac{1}{m} \sum_{i \in \bar{\mathcal{S}}} |\langle \Delta\theta, x_i \rangle| + \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle \Delta\theta, x_i \rangle - \varepsilon_i|,$$

Consider the following feasibility problem:

$$\text{find } \alpha \quad \text{s.t.} \quad \langle \alpha, x_i \rangle = 0, \forall i \in \bar{\mathcal{S}}, \quad \langle \alpha, x_i \rangle = \varepsilon_i, \forall i \in \mathcal{S}.$$

Since $m \leq d$ and $\{x_i\}_{i=1}^m$ are i.i.d. standard Gaussian vectors, they are linearly independent almost surely. Moreover, with high probability, at least one of ε_i 's will be nonzero. Therefore, with high probability, the above system of linear equations has at least one nonzero feasible solution. Suppose that $\bar{\alpha} \neq 0$ is one such solution. Define $\Delta\theta = \gamma \bar{\alpha} / \|\bar{\alpha}\|$ for some $0 < \gamma < \|\bar{\alpha}\|$. One can write

$$\mathcal{L}(\theta^\gamma + \Delta\theta) = \frac{1}{m} \left[1 - \frac{\gamma}{\|\bar{\alpha}\|} \sum_{i \in \mathcal{S}} |\varepsilon_i| \right] < \mathcal{L}(\theta^\gamma),$$

implying that $\Delta\theta$ is indeed a descent direction. Now, consider a 2-layer model. It is easy to verify the existence of a true solution $\mathbf{w} = (w_1, w_2)$ such that $w_1 \odot w_2 = \theta^\gamma$ and $\|w_1\|_0 = d$. Consider a perturbation of the form $\Delta\mathbf{w} = (0, \Delta w_2)$. One can write

$$\mathcal{L}(\mathbf{w} + \Delta\mathbf{w}) = \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle w_1 \odot \Delta w_2, x_i \rangle| + \frac{1}{m} \sum_{i \in \bar{\mathcal{S}}} |\langle w_1 \odot \Delta w_2, x_i \rangle - \varepsilon_i|.$$

Since w_1 is devoid of zero elements, there exists a nonzero Δw_2 such that $w_1 \odot \Delta w_2 = \gamma \bar{\alpha} / \|\bar{\alpha}\|$ for $\gamma < \|\bar{\alpha}\|$. Arguments analogous to 1-layer model can then be invoked to show that the constructed perturbation is indeed a descent direction. A similar idea can be naturally extended to $N \geq 3$.

Theorem 1 implies that, despite their convexity, 1-layer models are *not* suitable for the robust linear regression since the set of true solutions (which is a singleton $\mathcal{W} = \{\theta^?\}$) is unidentifiable. However, despite the existence of unidentifiable true solutions in N -layer models with $N \geq 2$, we will show that a simple SubGM converges to a *balanced* true solution, even if an arbitrarily large fraction of the measurements are corrupted with arbitrarily large noise values. This further sheds light on the desirable landscape of deeper models in the context of linear regression.

Algorithm 1 Sub-gradient Method

Input: Data points $\{(x_i, y_i)\}_{i=1}^m$, number of iterations T , the initial point \mathbf{w}_0 , and the step-size $\{\eta^{(t)}\}_{t=0}^T$;

Output: Solution $\mathbf{w}^{(T)}$ to (1);

for $t \leq T$ **do**

 Select a direction $\mathbf{d}^{(t)}$ from the sub-differential $\partial \mathcal{L}(\mathbf{w}^{(t)})$ defined as:

$$\partial_{w_i} \mathcal{L}(\mathbf{w}) = \frac{1}{m} \sum_{j=1}^n \text{Sign} \left(y_j - \sum_{k \neq i} w_k x_j \odot w_k \right); \quad (4)$$

 Set $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta^{(t)} \mathbf{d}^{(t)}$;

end for

3.2 Convergence of Sub-gradient Method

At every iteration t , SubGM selects a direction $\mathbf{d}^{(t)}$ from the sub-differential of the ℓ_1 -loss (defined as (4)), and updates the solution as $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta^{(t)} \mathbf{d}^{(t)}$; see Algorithm 1 for details. Our next two theorems characterize the performance of SubGM with small initialization on N -layer models. We consider the cases $N = 2$ and $N \geq 3$ separately, as SubGM behaves differently on these models. We define $\kappa = \theta_{\max}^? / \theta_{\min}^?$ as the condition number, where $\theta_{\max}^?$ and $\theta_{\min}^?$ are the maximum and minimum nonzero elements of $\theta^?$, respectively.

Theorem 2 (2-layer model). *Consider the iterations of SubGM $\{\mathbf{w}^{(t)}\}_{t=0}^T$ applied to $\mathcal{L}(\mathbf{w})$ with $N = 2$ and step-size $\eta \leq 1$. Suppose that the initial point satisfies $w_j^{(0)} = \Theta(\sqrt{\alpha} \mathbf{1})$, $j = 1, 2$, where $0 < \alpha \leq d^2 m / k$. Moreover, suppose that $m \geq \frac{k^2 - 2 \log^2(m) \log(d) \log(k^? k =)}{(1-p)^2}$. Then, the following statements hold with probability of $1 - Ce^{-\tilde{\Omega}(k)}$:*

- **Convergence guarantee:** After $\frac{1}{\eta} \log \frac{1}{\epsilon} \leq \bar{T} \leq \frac{k^3 - 2}{\eta} \log \frac{1}{\epsilon}$ iterations, we have

$$w_1^{(\bar{T})} \odot w_2^{(\bar{T})} - \theta^? \leq \eta \theta_{\max}^? \sqrt{d^2 m \alpha} \epsilon^{0.5} e^{-\frac{k^2}{(1-p)^2 m}}.$$

- **Balanced property:** For every $0 \leq t \leq \bar{T}$, we have

$$w_1^{(t)} - w_2^{(t)} \leq \frac{0.5}{\eta} \alpha \epsilon^{0.5} e^{-\frac{k^2}{(1-p)^2 m}}.$$

- **Long escape time:** For every $\bar{T} \leq t \leq \frac{1}{\eta} \frac{m(1-p)^2}{k} \bar{T}$, we have

$$w_1^{(t)} \odot w_2^{(t)} - \theta^? \leq \eta \theta_{\max}^? \sqrt{d^2 m \alpha} \epsilon^{0.5} e^{-\frac{k^2}{(1-p)^2 m}}.$$

Furthermore, if $m \geq d \log(m) / (1-p)^2$, with probability of $1 - Ce^{-\tilde{\Omega}(k)}$ and for every $t \geq \bar{T}$, we have

$$w_1^{(t)} \odot w_2^{(t)} - \theta^? \leq \eta \theta_{\max}^? \sqrt{d^2 m \alpha} \epsilon^{0.5} e^{-\frac{k^2}{m(1-p)^2}} (1 - \Omega(\eta \sqrt{d}))^t.$$

We provide the main idea behind the proof of Theorem 2 in Section 4. The formal proof can be found in the appendix. A few observations are in order based on Theorem 2. First, for any $\varepsilon > 0$, SubGM is guaranteed to satisfy $w_1^{(t)} \odot w_2^{(t)} - \theta^2 \leq \varepsilon$ after $\mathcal{O}((1/\varepsilon) \log(d/\varepsilon))$ iterations, provided that $\eta = \Theta(\varepsilon)$ and $\alpha = \varepsilon^2/(d^2m)$. Based on our numerical results (provided in the appendix), we believe that it is possible to establish a linear convergence for SubGM with a geometric step-size; a rigorous verification of this conjecture is considered future work. Second, although SubGM converges to a vicinity of a true solution quickly, it will stay there for a significantly longer time—in particular, $m(1-p)^2/k$ times longer than its initial convergence time. Such behavior is also exemplified in our simulations (see Figure 1e). After this escape time, the algorithm may slowly converge to an *overfitted* solution with a better training loss. Moreover, if $m \ll d$, SubGM will continuously converge to a true solution at an exponential rate, and it will never diverge. Finally, Theorem 2 shows that SubGM implicitly favors balanced solutions, i.e. solutions whose factors have similar magnitudes. Combined with the convergence result of SubGM, we immediately conclude that SubGM converges to a particular solution of the form $(\sqrt{\theta^2}, \sqrt{\theta^2})$. Therefore, the solution found by SubGM will enjoy the same (approximate) sparsity pattern as θ^2 .

Theorem 3 (N -layer models). *Consider the iterations of SubGM $\{\mathbf{w}^{(t)}\}_{t=0}^T$ applied to $\mathcal{L}(\mathbf{w})$ with $N \geq 3$ and step-size $\eta \leq \frac{1}{N} \alpha \frac{N-2}{N}$. Suppose that the initial point satisfies $w_j^{(0)} = \Theta(\alpha^{1-N} \mathbf{1})$, where $0 < \alpha \leq d^2m/k$. Moreover, suppose that $m \ll \frac{k^2 \log^2(m) \log(d) \log(k^2/k)}{(1-p)^2}$. Then, the following statements hold with probability of $1 - Ce^{-\tilde{\Omega}(k)}$:*

- **Convergence guarantee:** After $\frac{1}{N} \alpha \frac{N-2}{N} \cdot \bar{T} \cdot \frac{k^{3-2}}{N} \alpha \frac{N-2}{N}$ iterations, we have

$$\forall w_i^{(\bar{T})} - \theta^2 \leq N\eta\theta_{\max}^2 \vee \sqrt{d^2m\alpha}.$$

- **Balanced property:** For every $0 \leq t \leq \bar{T}$, we have

$$\begin{aligned} w_{i;l}^{(t)} - w_{j;l}^{(t)} &= \mathcal{O}(\alpha^{1-N}), & \text{for } 1 \leq i < j \leq N, l : \theta_l^2 = 0, \\ w_{i;l}^{(t)} - w_{j;l}^{(t)} &= \tilde{\mathcal{O}}\left(\frac{\theta_l^2}{k^3/m}\right), & \text{for } 1 \leq i < j \leq N, l : \theta_l^2 \neq 0. \end{aligned}$$

- **Long escape time:** For every $\bar{T} \leq t \leq \frac{1}{k} \frac{m(1-p)^2}{k} \bar{T}$, we have

$$\forall w_i^{(t)} - \theta^2 \leq N\eta\theta_{\max}^2 \vee \sqrt{d^2m\alpha},$$

Furthermore, if $m \ll d^{\frac{2N-2}{N}} \log(m)/(1-p)^2$, with probability of at least $1 - Ce^{-\tilde{\Omega}(k)}$ and for every $t > \bar{T}$, we have

$$\forall w_i^{(t)} - \theta^2 \leq N\eta\theta_{\max}^2 \vee \frac{\sqrt{d^2m\alpha}}{\sqrt{d^2m\alpha N\eta d^{\frac{N-1}{N}}(t-\bar{T})+1}} \frac{1}{N^{\frac{N-2}{N}}}.$$

The proof of this theorem can be found in the appendix. Theorem 3 sheds light on an important benefit of N -layer models with $N \geq 3$ compared to 2-layer models: for sufficiently small step-size, deeper models improve the generalization error by a factor of $(1/\alpha)^{\tilde{\Theta}(k^2/\sqrt{(1-p)^2m})}$. This improvement is particularly significant when both α and m are small. However, such improvement comes at the expense of a slower convergence rate. In particular, after setting $\eta = \Theta(\varepsilon/N)$, and $\alpha = \varepsilon/\sqrt{d^2m}$, SubGM needs $\mathcal{O}((1/\varepsilon)^{1+\frac{N-2}{N}})$ iterations to obtain an ε -accurate solution. Evidently, the convergence rate deteriorates with N , ultimately approaching $\mathcal{O}(1/\varepsilon^2)$ for infinitely deep models. This can be observed in practice: Figures 1e and 1f show that 3-layer model enjoys a better generalization error compared to 2-layer model, but suffers from a slower convergence rate. This slower convergence rate also manifests itself in a more stable behavior of the algorithm: for deeper models, SubGM stays close to the ground truth for a longer time. Finally, the balanced property of the solution obtained via SubGM extends to N -layer models. In particular, SubGM converges to a particular solution of the form $(\sqrt[2]{\theta^2}, \dots, \sqrt[2]{\theta^2})$, thereby inheriting the same sparsity pattern as θ^2 .

3.3 Local Landscape Around Balanced Solution

In the previous section, we showed that SubGM converges to a balanced solution. In this section, we characterize the local landscape around this balanced solution, proving that it becomes flatter for deeper models.

Theorem 4 (flatness around balanced solution). *Suppose that $k \log(d)/(1-2p)^2 \cdot m \leq 0.1d$ and $p < 1/2$. Let $\mathbf{w}^\gamma = (\sqrt[\gamma]{\theta^\gamma}, \dots, \sqrt[\gamma]{\theta^\gamma})$. Then, for any $N \geq 2$ and $\gamma \leq t_0/\sqrt{d} \wedge 1$, the following statements hold:*

- With probability at least $1 - e^{-\Omega(k)}$, we have

$$\inf_{\mathbf{w}:k\mathbf{w}} \inf_{\mathbf{w}^\gamma k_\gamma} \{\mathcal{L}(\mathbf{w}^\gamma) - \mathcal{L}(\mathbf{w})\} \leq -\frac{d}{\sqrt{m}} \gamma^N.$$

- With probability at least $1/16$, we have

$$\inf_{\mathbf{w}:k\mathbf{w}} \inf_{\mathbf{w}^\gamma k_\gamma} \{\mathcal{L}(\mathbf{w}^\gamma) - \mathcal{L}(\mathbf{w})\} \geq -\sqrt{p_0 p} \frac{d}{\sqrt{m}} \gamma^N.$$

Theorem 4 shows that, within a γ -neighborhood of \mathbf{w}^γ , the most descent direction from \mathbf{w}^γ can reduce the loss by at most $\mathcal{O}(d/\sqrt{m} \cdot \gamma^N)$, which decreases exponentially with N . Moreover, in the noisy setting, the above theorem implies that \mathbf{w}^γ is likely to be neither local nor global minimum since it has a descent direction. However, the flatness of the landscape around \mathbf{w}^γ enables SubGM to stay close to the balanced solution for a long time.

Remark 1. *Note that the choice of ℓ_γ -ball for the perturbation set is to ensure that the size of the possible perturbations per layer remains independent of the depth of the model. This is indeed crucial to ensure a fair comparison between models with different depths: alternative choices of the perturbation set, such as ℓ_q -ball with $1 \leq q < \infty$ (e.g. ℓ_2 -ball) would shrink the size of the feasible per-layer perturbations with N , thereby leading to an unfair advantage to deeper models.*

4 Proof Techniques

At the crux of our proof technique for Theorems 2 and 3 lies the following decomposition of the sub-differential:

$$\partial \mathcal{L}(\mathbf{w}) = \underbrace{\xi \cdot \partial \bar{\mathcal{L}}(\mathbf{w})}_{\text{expected subdiff.}} + \underbrace{\partial \mathcal{L}(\mathbf{w}) - \xi \cdot \partial \bar{\mathcal{L}}(\mathbf{w})}_{\text{subdiff. deviation}}, \quad \text{for some strictly positive } \xi.$$

In the above decomposition, $\bar{\mathcal{L}}(\mathbf{w})$ is called *expected loss*, and is defined as $\bar{\mathcal{L}}(\mathbf{w}) = \|w_1 \odot \dots \odot w_N - \theta^\gamma\|$. As will be shown later, $\bar{\mathcal{L}}(\mathbf{w})$ captures the expected behavior of the empirical loss $\mathcal{L}(\mathbf{w})$. To analyze the behavior of SubGM on $\mathcal{L}(\mathbf{w})$, we first consider the ideal scenario, where $\mathcal{L}(\mathbf{w})$ coincides with its expectation. Then, we extend our analysis to the general case by controlling the sub-differential deviation. In particular, we show that the desirable convergence properties of SubGM extends to $\mathcal{L}(\mathbf{w})$, provided that its sub-differentials are “direction-preserving”, i.e., $\mathbf{d} \approx \xi \bar{\mathbf{d}}$, for every $\mathbf{d} \in \partial \mathcal{L}(\mathbf{w})$, $\bar{\mathbf{d}} \in \partial \bar{\mathcal{L}}(\mathbf{w})$ and some $\xi > 0$. To formalize this idea, we first provide a more concise characterization of $\partial \mathcal{L}(\mathbf{w})$:

$$\partial_{w_i} \mathcal{L}(\mathbf{w}) = \left\{ \begin{array}{c} \underbrace{q \odot w_k}_{k \neq i} : q \in \mathcal{Q} \\ \underbrace{\theta^\gamma - w_k}_{k} \end{array} \right\}, \quad \text{where } \mathcal{Q}(z) = \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, z \rangle + \varepsilon_i) x_i.$$

Definition 1 (approximately sparse vectors). *We say a vector $v \in \mathbb{R}^d$ is (k, ϑ) -approximately sparse if there exists a vector u , such that $\|u\|_0 \leq k$ and $\|u - v\| \leq \vartheta$.*

Proposition 1 (direction-preserving property). *Suppose that $m \geq \frac{k \log^2(m) \log(d) \log(R/\delta)}{(1-p)^2}$ for some $R, \vartheta, \delta > 0$. Then, with probability of at least $1 - Ce^{-\Omega(m^2)}$, the following inequality holds for any $q \in \mathcal{Q}(z)$ and any (k, ϑ) -approximately sparse vector z that satisfies $\frac{dm}{k\vartheta} \log(1/\vartheta) \cdot \|z\| \leq R$:*

$$q - \frac{1}{\pi} \left(1 - p + p e^{-\frac{2k\vartheta}{\|z\|}} \right) \leq \delta. \quad (5)$$

Moreover, if $m \geq \frac{d \log(m)}{(1-p)^2}$, with probability of $1 - Ce^{-\Omega(m^2)}$, (5) holds for every $z \in \mathbb{R}^d$.

Proposition 1 is analogous to *Sign-RIP* condition introduced in [29, 28] for the robust low-rank matrix recovery, and is at the heart of our proofs for Theorems 2 and 3. Suppose that $\theta_1^? \dots \theta_k^?$ is a (k, ϑ) -approximately sparse and satisfies (5). Then, we have $\mathbf{d} - \bar{\mathbf{d}}_1 \leq \max_{i \neq k} w_k \delta$, which in turn provides an upper bound on the sub-differential deviation.

4.1 Proof Sketch of Theorem 2

To streamline the presentation, here we only provide simplified versions of our key ideas, which inevitably lead to looser guarantees. To streamline the proof, we assume that $\theta_1^? \geq \dots \geq \theta_k^? > \theta_{k+1}^? = \dots = \theta_d^? = 0$. Moreover, for simplicity of notation, we denote $u = w_1$ and $v = w_2$. Consider the following decomposition:

$$u \odot v = \left[\underbrace{u_1 v_1 \dots u_k v_k}_S \mid \underbrace{u_{k+1} v_{k+1} \dots u_d v_d}_E \right]^>. \quad (6)$$

The vectors S and E are called *signal* and *residual terms*, respectively. Evidently, we have $u \odot v = \theta^?$ if and only if $S = [\theta_1^?, \dots, \theta_k^?]^>$ and $E = 0$. Based on this observation, our goal is to show that the signal term converges to $[\theta_1^?, \dots, \theta_k^?]^>$ exponentially fast, while the error term remains small throughout the solution trajectory.

Lemma 1 (signal dynamic; informal). *Suppose that (5) holds for $z = \theta^? - u^{(t)} \odot v^{(t)}$, and $\theta^? - u^{(t)} \odot v^{(t)} \ll \eta \|\theta^?\|$. Then, we have*

$$u_i^{(t+1)} v_i^{(t+1)} \geq 1 + 2\eta \frac{\theta_i^? - u_i^{(t)} v_i^{(t)}}{u^{(t)} \odot v^{(t)} - \theta^?} + \delta_i \quad u_i^{(t)} v_i^{(t)}, \quad (7)$$

for some $|\delta_i| \leq \delta$ and every $i = 1, \dots, k$.

Lemma 2 (residual dynamic; informal). *Suppose that (5) holds for $z = \theta^? - u^{(t)} \odot v^{(t)}$, and $\theta^? - u^{(t)} \odot v^{(t)} \ll \eta \|\theta^?\|$. Then, we have*

$$u_i^{(t+1)}{}^2 + v_i^{(t+1)}{}^2 \leq (1 + \mathcal{O}(\eta\delta)) \quad u_i^{(t)}{}^2 + v_i^{(t)}{}^2, \quad (8)$$

for every $i = k + 1, \dots, d$.

Lemma 3 (difference dynamic; informal). *Suppose that (5) holds for $z = \theta^? - u^{(t)} \odot v^{(t)}$, and $\theta^? - u^{(t)} \odot v^{(t)} \ll \eta \|\theta^?\|$. Then, we have*

$$u_i^{(t+1)} - v_i^{(t+1)} = u_i^{(t)} - v_i^{(t)} \left(1 - \eta \frac{\theta_i^? - u_i^{(t)} v_i^{(t)}}{u^{(t)} \odot v^{(t)} - \theta^?} + \eta \delta_i \right), \quad (9)$$

for some $|\delta_i| \leq \delta$ and every $i = 1, \dots, d$.

Convergence guarantee. For any fixed $i = 1, \dots, k$, we show that $u_i^{(t)} v_i^{(t)} = \theta_i^? \pm \mathcal{O}(\delta) \|\theta^?\|$ after $\mathcal{O}(\|\theta^?\|/(\eta\theta_i^?) \log(1/\alpha))$ iterations. To see this, suppose that T_i is the largest iteration such that $u_i^{(t)} v_i^{(t)} \leq \theta_i^?$ for every $t \leq T_i$. Moreover, suppose that $u^{(t)} \odot v^{(t)} \leq C \|\theta^?\|$, for sufficiently large C (this is proven in the appendix). Under these assumptions, (7) reduces to

$$u_i^{(t+1)} v_i^{(t+1)} \geq 1 + \Omega(1) \frac{\eta\theta_i^?}{\|\theta^?\|} \quad u_i^{(t)} v_i^{(t)}. \quad (10)$$

which implies that $T_i \leq \|\theta^?\|/(\eta\theta_i^?) \log(1/\alpha)$. For any $t > T_i$, define $y_i^{(t)} = \theta_i^? - u_i^{(t)} v_i^{(t)}$. One can write

$$y_i^{(t+1)} \leq 1 - \Omega(1) \frac{\eta\theta_i^?}{\|\theta^?\|} \quad y_i^{(t)} + \eta\delta\theta_i^?. \quad (11)$$

Hence, with additional $\mathcal{O}(\|\theta^?\|/(\eta\theta_1^?))$ iterations, we have $u_i^{(t)} v_i^{(t)} = \theta_i^? \pm \mathcal{O}(\delta) \|\theta^?\|$. On the other hand, Lemma 2 implies that, for any $i = k + 1, \dots, d$ and $t \leq \|\theta^?\|/(\eta\theta_k^?) \log(1/\alpha)$, we have

$$u_i^{(t)}{}^2 + v_i^{(t)}{}^2 \leq \alpha (1 + \mathcal{O}(\eta\delta))^{\mathcal{O}(\frac{\|\theta^?\|}{\eta\theta_k^?} \log(1/\alpha))} \leq \alpha^1 \mathcal{O}(\alpha^{\frac{1}{k}}),$$

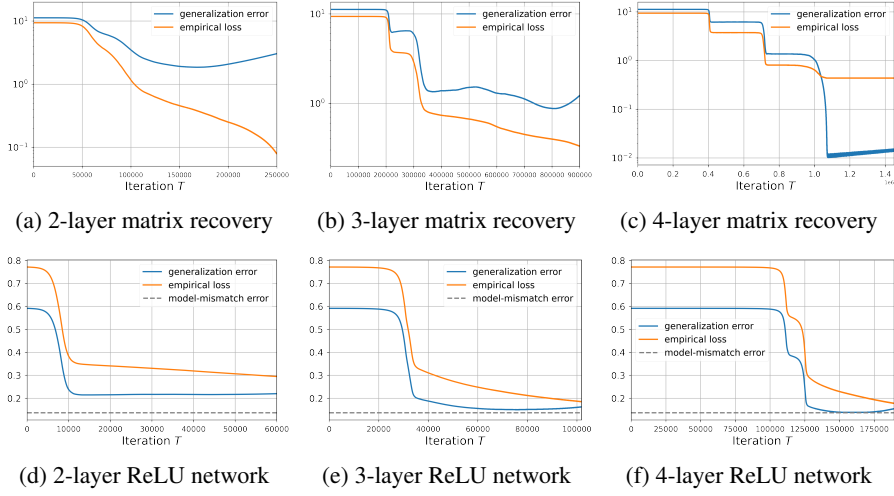


Figure 2: Deep matrix recovery (first row). The ground truth $X^? \in \mathbb{R}^{20 \times 20}$ with $\text{rank}(X^?) = 3$ is chosen randomly. The elements of the measurement matrices are selected from $\mathcal{N}(0; 1)$, and the sample size is set to $m = 180$. The corruption probability is set to $p = 0.05$ with distribution $\mathcal{N}(0; 100)$. We use SubGM with step-size $\eta = 0.001$ and Gaussian initialization with an initialization scale $\sigma = 1 \cdot 10^{-3}$. **ReLU models (second row).** The samples are chosen from $y_i = \sin(\theta^? x_i) + u_i$ where $\theta^? \in \mathbb{R}^{50}$ is randomly generated with $k^? k_0 = 2$, $x_i \sim \mathcal{N}(0; I_{50})$, and $u_i \sim \mathcal{N}(0; 25)$ with corruption probability $p = 0.05$. The sample size is set to be $m = 1500$. We use SubGM with step-size $\eta = 0.001$ and Gaussian initialization $\mathcal{N}(0; \sigma^2 = d)$.

where $\kappa = \theta_1^? / \theta_k^?$ is the condition number of $\theta^?$. Combining the above dynamics, we have

$$u^{(t)} \odot v^{(t)} - \theta^? \leq \eta \|\theta^?\| \sqrt{k} \|\theta^?\| \delta \sqrt{d} \alpha^{1 - \mathcal{O}(\frac{P}{k})}.$$

In the appendix, we provide a more refined analysis that relaxes the dependency of the final error on δ and κ .

Long escape time. We show in the appendix that after the first stage, the residual becomes the dominant term in the final error. This together with Lemma 2 implies that, for every $t \geq \frac{k^? k}{\eta \delta} \log(1/\alpha)$, we have $\|E\| \leq \sqrt{d} \alpha^{1 - \mathcal{O}(\frac{P}{k})}$.

Balanced property. We have $u_i^{(t)} v_i^{(t)} \leq \theta_i^?$ for every $i \in [k]$, and $|u_i^{(t)} v_i^{(t)}| \leq \alpha^{1 - \mathcal{O}(\frac{P}{k})}$ for every $i = k + 1, \dots, d$. Therefore, Lemma 3 can be invoked to verify $|u_i^{(t+1)} - v_i^{(t+1)}| \leq (1 + \mathcal{O}(\eta \delta)) |u_i^{(t)} - v_i^{(t)}|$. This in turn leads to

$$|u_i^{(t)} - v_i^{(t)}| \leq \sqrt{\alpha} (1 + \mathcal{O}(\eta \delta)) \alpha^{\frac{k^? k}{k} \log(1/\alpha)} \leq \alpha^{0.5} \alpha^{\mathcal{O}(\frac{P}{k})}.$$

5 Numerical Experiments: Beyond Linear Regression

In this section, we empirically verify that the benefits of depth extend to the robust variants of deep matrix recovery and ReLU networks with ℓ_1 -loss.

Deep Matrix Recovery. In low-rank matrix recovery, the goal is to recover a low-rank matrix $X^? \in \mathbb{R}^{d \times d}$, from a limited number of noisy measurements of the form $y_i = \langle A_i, X^? \rangle + \varepsilon_i$. To recover $X^?$, we consider a deep factorized model of the form $W_1 W_2 \dots W_N$, where $W_i \in \mathbb{R}^{d \times d}$ for $i = 1, \dots, N$, and minimize the ℓ_1 -loss $(1/m) \sum_{i=1}^m |y_i - \langle A_i, W_1 W_2 \dots W_N \rangle|$ via SubGM. When $N = 2$, the above model reduces to the famous Burer-Monteiro approach [8]. We assume that 5% of the measurements are grossly corrupted with noise. The first row of Figure 2 shows the performance of SubGM on 2-, 3-, and 4-layer models. It can be seen that the 4-layer model outperforms shallower models, achieving a generalization error that is proportional to the step-size.

Deep ReLU Network on Synthetic Dataset. As another experiment, we analyze the effect of depth on the performance of SubGM with ReLU networks and ℓ_1 -loss. Given an input $x \in \mathbb{R}^d$, the output of an N -layer ReLU network is defined as $f_{\mathbf{W}}(x) = W_N \sigma(W_{N-1} \cdots \sigma(W_1 x) \cdots)$, where $W_1 \in \mathbb{R}^{m \times d}, W_2, \dots, W_{N-1} \in \mathbb{R}^{m \times m}$, and $W_N \in \mathbb{R}^{1 \times m}$. Moreover, $\sigma(x) = \max\{0, x\}$ is the ReLU activation function. Given the true function $f^?(x) = \sin(\theta^? x)$, our goal is to train a ReLU model to approximate $f^?$ as accurately as possible. To this goal, we minimize the ℓ_1 -loss $(1/m) \sum_{i=1}^m |y_i - f_{\mathbf{W}}(x_i)|$. The second row of Figure 2 illustrates the performance of SubGM. It is worth noting that, unlike robust linear regression and deep matrix recovery, there always exists a non-diminishing model-mismatch error between the true and considered ReLU model (shown as a dashed line). Nonetheless, SubGM can achieve this model-mismatch error on a 4-layer ReLU model with only 1500 samples, even if 5% of the measurements are corrupted with large noise.

Deep ReLU Network on CIFAR Dataset. We verify that the desirable performance of SubGM with ℓ_1 -loss can be extended to its stochastic variant with mini-batches on CIFAR-10 and CIFAR-100 [24], outperforming cross-entropy (CE) loss, which is considered as one of the most suitable loss functions for CIFAR datasets. To show this, we use standard ResNet architectures [21] with ℓ_1 -loss and compare it with the cross-entropy loss on noisy CIFAR datasets, where we randomize the labels of 10% of the training dataset. For CIFAR-100 experiment, we use the ‘‘loss scaling’’ trick introduced in [22]. The training details are deferred to Section B.3. The best test accuracy for both CIFAR-10 and CIFAR-100 is reported in Table 1. One can see that ℓ_1 -loss outperforms cross-entropy loss significantly, demonstrating that our framework may be extended to more realistic settings. Moreover, we do observe that the deeper model performs better on CIFAR-100, which aligns with our theoretical result. Based on our simulations, an interesting and important future direction would be to study the optimization landscape of ℓ_1 -loss with more general neural network architectures.

	CIFAR-10			CIFAR-100		
	ResNet-18	ResNet-34	ResNet-50	ResNet-18	ResNet-34	ResNet-50
CE loss	91.52%	91.53%	90.87%	70.17%	71.22%	71.30%
ℓ_1 -loss	94.16%	93.13%	92.68%	73.14%	74.86%	76.46%

Table 1: Test accuracy for ResNets on CIFAR-10 and CIFAR-100 datasets with 10% label noise.

6 Conclusion

Modern problems in machine learning are naturally nonconvex but can be solved reasonably well in practice. To explain this, a recent body of work has postulated that many optimization problems in machine learning are ‘‘convex-like’’, i.e., they are devoid of spurious local minima. Our work shows that such global property is too restrictive to hold even in the context of linear regression, and instead propose a more refined *trajectory analysis* to better capture the landscape of the problem around the solution trajectory. We show that convex models may be fundamentally ill-suited for linear models, and deeper models—despite their nonconvexity—have provably better optimization landscape around the solution trajectory. Empirically, we show that our analysis may extend beyond linear regression; formal verification of this conjecture is considered an enticing challenge for future research.

Acknowledgements

We thank Richard Y. Zhang, Cédric Józsz, and Tiffany Wu for helpful discussions and feedback. We thank Ruiqi Gao and Jiaye Teng for insightful discussions in the initial phase of this work. We are also thankful to an anonymous reviewer for pointing out the relationship between the perturbation ball and the depth of linear models. This research is supported, in part, by NSF Award DMS-2152776, ONR Award N00014-22-1-2127, MICDE Catalyst Grant, MIDAS PODS grant, and Startup funding from the University of Michigan.

References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- [2] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32:7413–7424, 2019.
- [3] Richard Baraniuk, Mark Davenport, Ronald DeVore, and Michael Wakin. A simple proof of the restricted isometry property for random matrices. *Constructive Approximation*, 28(3):253–263, 2008.
- [4] Dimitris Bertsimas, Angela King, and Rahul Mazumder. Best subset selection via a modern optimization lens. *The annals of statistics*, 44(2):813–852, 2016.
- [5] Srinadh Bhojanapalli, Behnam Neyshabur, and Nathan Srebro. Global optimality of local search for low rank matrix recovery. *arXiv preprint arXiv:1605.07221*, 2016.
- [6] Peter J Bickel, Ya’acov Ritov, and Alexandre B Tsybakov. Simultaneous analysis of lasso and dantzig selector. *The Annals of statistics*, 37(4):1705–1732, 2009.
- [7] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [8] Samuel Burer and Renato DC Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Mathematical Programming*, 95(2):329–357, 2003.
- [9] Emmanuel Candes and Terence Tao. The dantzig selector: Statistical estimation when p is much larger than n . *The annals of Statistics*, 35(6):2313–2351, 2007.
- [10] Hung-Hsu Chou, Johannes Maly, and Holger Rauhut. More is less: Inducing sparsity via overparameterization. *arXiv preprint arXiv:2112.11027*, 2021.
- [11] Damek Davis and Dmitriy Drusvyatskiy. Stochastic subgradient method converges at the rate $o(k^{-1/4})$ on weakly convex functions. *arXiv preprint arXiv:1802.02988*, 2018.
- [12] Lijun Ding, Liwei Jiang, Yudong Chen, Qing Qu, and Zhihui Zhu. Rank overspecified robust matrix recovery: Subgradient method and exact recovery. *Advances in Neural Information Processing Systems*, 34, 2021.
- [13] Simon Du and Wei Hu. Width provably matters in optimization for deep linear neural networks. In *International Conference on Machine Learning*, pages 1655–1664. PMLR, 2019.
- [14] Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. *arXiv preprint arXiv:1806.00900*, 2018.
- [15] Ronen Eldan and Ohad Shamir. The power of depth for feedforward neural networks. In *Conference on learning theory*, pages 907–940. PMLR, 2016.
- [16] Salar Fattahi and Somayeh Sojoudi. Exact guarantees on the absence of spurious local minima for non-negative rank-1 robust principal component analysis. *Journal of machine learning research*, 2020.
- [17] Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [18] Daniel Gissin, Shai Shalev-Shwartz, and Amit Daniely. The implicit bias of depth: How incremental learning drives generalization. *arXiv preprint arXiv:1909.12051*, 2019.
- [19] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Implicit bias of gradient descent on linear convolutional networks. *arXiv preprint arXiv:1806.00468*, 2018.

- [20] Jeff Z HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. Shape matters: Understanding the implicit bias of the noise covariance. In *Conference on Learning Theory*, pages 2315–2357. PMLR, 2021.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [22] Like Hui and Mikhail Belkin. Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks. In *International Conference on Learning Representations*, 2020.
- [23] Kenji Kawaguchi. Deep learning without poor local minima. *arXiv preprint arXiv:1605.07110*, 2016.
- [24] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [25] Jianguan Li, Thanh Nguyen, Chinmay Hegde, and Ka Wai Wong. Implicit sparse regularization: The impact of depth and early stopping. *Advances in Neural Information Processing Systems*, 34, 2021.
- [26] Xiao Li, Zhihui Zhu, Anthony Man-Cho So, and Rene Vidal. Nonconvex robust low-rank matrix recovery. *SIAM Journal on Optimization*, 30(1):660–686, 2020.
- [27] Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. *arXiv preprint arXiv:2012.09839*, 2020.
- [28] Jianhao Ma and Salar Fattahi. Sign-rip: A robust restricted isometry property for low-rank matrix recovery. *arXiv preprint arXiv:2102.02969*, 2021.
- [29] Jianhao Ma and Salar Fattahi. Global convergence of sub-gradient method for robust matrix recovery: Small initialization, noisy measurements, and over-parameterization. *arXiv preprint arXiv:2202.08788*, 2022.
- [30] Rahul Mazumder, Peter Radchenko, and Antoine Dedieu. Subset selection with shrinkage: Sparse linear modeling when the snr is low. *arXiv preprint arXiv:1708.03288*, 2017.
- [31] Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Sensitivity and generalization in neural networks: an empirical study. *arXiv preprint arXiv:1802.08760*, 2018.
- [32] Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE transactions on information theory*, 57(10):6976–6994, 2011.
- [33] Itay Safran and Jason D Lee. Optimization-based separations for neural networks. *arXiv preprint arXiv:2112.02393*, 2021.
- [34] Matus Telgarsky. Representation benefits of deep feedforward networks. *arXiv preprint arXiv:1509.08101*, 2015.
- [35] Matus Telgarsky. Benefits of depth in neural networks. In *Conference on learning theory*, pages 1517–1539. PMLR, 2016.
- [36] Sara Van de Geer. The deterministic lasso. Seminar für Statistik, Eidgenössische Technische Hochschule (ETH) Zürich, 2007.
- [37] Aad W Van Der Vaart, Adrianus Willem van der Vaart, Aad van der Vaart, and Jon Wellner. *Weak convergence and empirical processes: with applications to statistics*. Springer Science & Business Media, 1996.
- [38] Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization for optimal sparse recovery. *Advances in Neural Information Processing Systems*, 32:2972–2983, 2019.

- [39] Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- [40] Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pages 3635–3673. PMLR, 2020.
- [41] Peng Zhao, Yun Yang, and Qiao-Chu He. Implicit regularization via hadamard product over-parametrization in high-dimensional linear regression. *arXiv preprint arXiv:1903.09367*, 2019.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#)
 - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#)
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) The complete proofs are provided in the appendix.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[No\]](#) We provided all the details for our experiments to be reproducible.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[N/A\]](#)
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) The details are provided in the appendix.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[N/A\]](#)
 - (b) Did you mention the license of the assets? [\[N/A\]](#)
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

Contents

A Related Works	14
B Additional Experiments	15
B.1 Deeper Linear Models	15
B.2 Geometric Step-size	15
B.3 Experiments on CIFAR Dataset	16
C Proofs of Landscape Analysis	17
C.1 Proof of Theorem 1	17
C.2 Proof of Theorem 4	19
D Proofs of Convergence Analysis	22
D.1 Proof of Theorem 2	22
D.2 Proof of Theorem 3	25
E Proof of Proposition 1	29
F Auxiliary Lemmas	34
G Deferred Proofs	35
G.1 Proof of Lemma 5	35
H Preliminaries on the Uniform Concentration Bounds	35

A Related Works

Deep models: It is known that deeper models enjoy better *approximation power*. For instance, [15, 34, 35] introduce several functions that are expressible by deep models of moderate size; yet, they cannot be approximated via any shallow network of sub-exponential size. A recent work [33] shows that depth separation may lead to optimization separation. In other words, functions that can be expressed by deeper models can also be efficiently learned via gradient descent. Another line of work shows that deep linear models have a strong implicit bias towards the true solution [19, 2, 10], and that they benefit from *incremental learning* [18, 27]. More generally, [1] show that stochastic gradient descent on deep nonlinear models can provably learn certain complex functions by automatically decomposing them into a series of simpler ones.

Robust and sparse linear regression: Robust and sparse linear regression is a classical problem in statistics, with a wide range of applications in signal and image processing. Regularized methods, including Lasso [36, 9, 6, 32], Best Subset Selection [30, 4], and Forward and Backward Step-wise Regression [17], are considered as most widely used methods for solving robust and sparse linear regression that come equipped with strong statistical and computational guarantees. Recently, sparse linear regression has been used to explore the implicit bias of different optimization algorithms and initialization regimes. [38, 41] show that, for some unregularized overparameterized models, gradient descent (with early stopping) for sparse linear regression achieves minimax error rate. Moreover, [40] study how the scale of the initial point controls the transition between the “kernel” (lazy training) and “rich” regimes, and their corresponding generalization performance. [20] use this problem setting to explore the role of label noise in stochastic gradient descent.

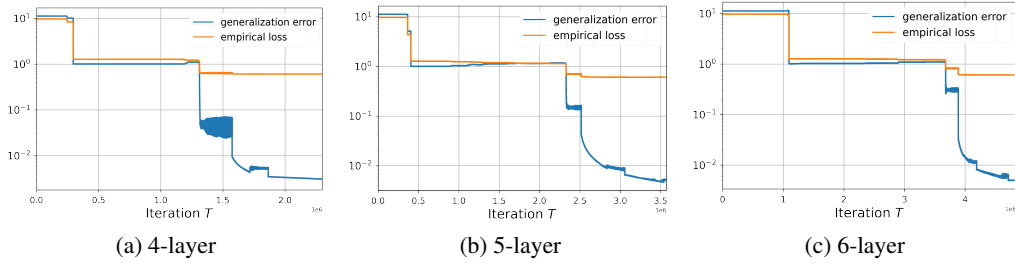


Figure 3: The optimization trajectories of deep models ($N = 4; 5; 6$).

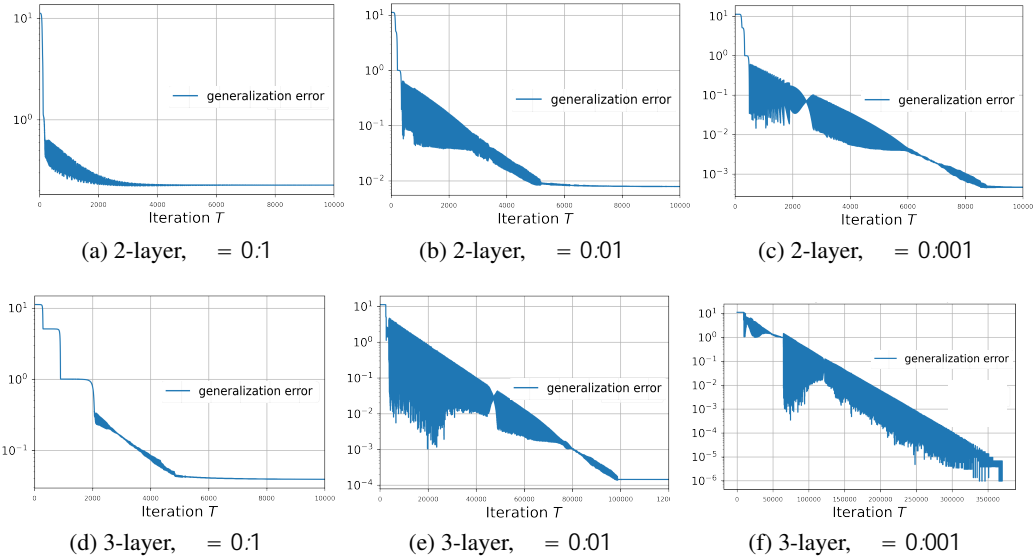


Figure 4: The optimization trajectories of 2 and 3-layer models with different initialization size $= 0.1; 0.01; 0.001$ using exponentially decayed step-size.

B Additional Experiments

In this section, we provide additional experiments on the performance of SubGM on deep models. Our goal is to verify our theoretical results and show the benefits of both small initialization and geometric step-size. Moreover, we show that the desirable performance of SubGM can be observed in its stochastic variant, as well as for different architectures of ResNets with ℓ_1 -loss, and more realistic CIFAR-10 dataset. All simulations are run on a desktop computer with an Intel Core i9 3.50 GHz CPU and 128GB RAM. The reported results are for implementation in Python.

B.1 Deeper Linear Models

In this experiment, we study the performance of SubGM on deeper models ($N = 4, 5, 6$). To accelerate the training process, we first use a large step-size $\eta_1 = 1 \times 10^{-3}$, and then progressively apply smaller step-sizes $\eta_2 = 1 \times 10^{-4}$ and $\eta_3 = 1 \times 10^{-5}$ as the training loss continues to decay. As shown in Figure 3, deeper models share similar generalization error, outperforming 1- and 2-layer models presented in Figure 1.

B.2 Geometric Step-size

As shown in our simulations, training N -layer models may require millions of iterations even on a small synthetic dataset. As proven in Theorem 3, the training process may become even slower for deeper models. In this experiment, we explore the performance of geometric (i.e., exponentially

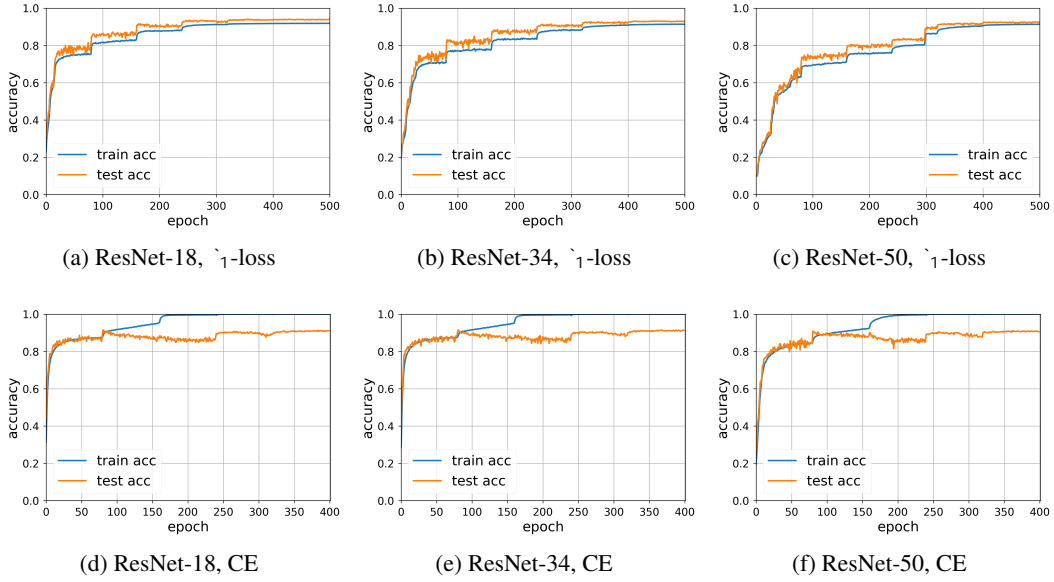


Figure 5: We apply ResNet-18, 34, 50 on noisy CIFAR-10 with both ℓ_1 -loss and cross-entropy loss (CE). For the training dataset, we randomly choose 10% samples and replace their labels with uniform random labels. We use SGD with initial learning rate $\eta = 0.1$, momentum 0.9, batch size $B = 32$. For every 80 epochs, we decay the learning rate by a factor 0.33. We use standard data augmentation. The initialization is set by default in PyTorch.

decaying) step-size on the same dataset. SubGM with a geometric step-size has been widely used for the optimization of ℓ_1 -loss [26, 28], and more general sharp weakly convex functions [11]. Figure 4 shows that a geometric step-size can lead to a 1000-fold reduction in the required number of iterations. Moreover, a geometric step-size improves the convergence rate to linear. The theoretical justification of this improvement is left as an enticing challenge for future research. Finally, it can be observed that SubGM with geometric step-size performs surprisingly well on deeper models, achieving a generalization error in the order of 10^{-6} . This further supports the benefits of depth.

B.3 Experiments on CIFAR Dataset

In this section, we provide the training details for the experiments on both CIFAR-10 and CIFAR-100 where 10% of the training data points are randomly labeled. For CIFAR-100 experiment, we use the “loss scaling” trick introduced in [22]. In particular, we denote the neural network by $f : \mathbb{R}^d \rightarrow \mathbb{R}^C$, where d is the input dimension and C is the number of class. The standard ℓ_1 -loss for the one-hot encoded label vector can be written (at a single point) as

$$\ell = |f(x)[c] - 1| + \sum_{c' \neq c} |f(x)[c']|. \quad (12)$$

Here c is the position of the label and $f(x)[i]$ is the i -th coordinate of the prediction. The rescaled ℓ_1 -loss is defined by two parameters k and M as follows:

$$\ell_{\text{scaling}} = k \cdot |f(x)[c] - M| + \sum_{c' \neq c} |f(x)[c']|. \quad (13)$$

In our simulation, we choose $k = 5$, and $M = 2$. Moreover, in the CIFAR-100 experiment, we use ℓ_2 -loss with the same k, M in the first 30 epochs as the warm-up phase. The evolution of the training and test accuracy for CIFAR-10 and CIFAR-100 with both ℓ_1 and CE losses are shown in Figures 5 and 6.

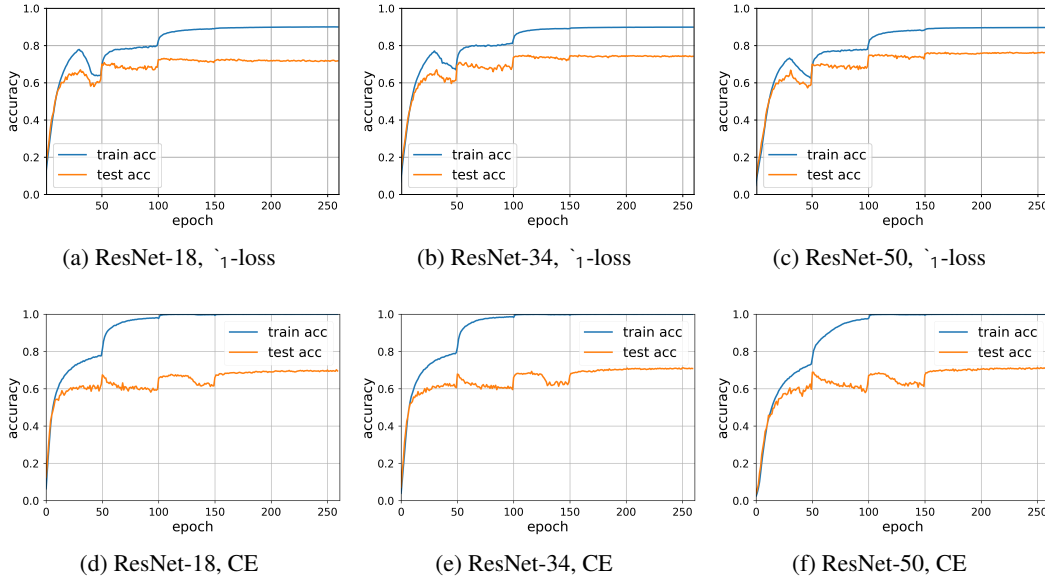


Figure 6: We apply ResNet-18, 34, 50 on noisy CIFAR-100 with both ℓ_1 -loss and cross-entropy loss (CE). The data generator and optimizer are the same as those in CIFAR-10 experiment. For each trial, we set the initial learning rate to be $\eta = 0.05$ and run 260 epochs and decay the learning rate with factor 0.33 for each 50 epochs.

C Proofs of Landscape Analysis

C.1 Proof of Theorem 1

Over-parameterized Regime

We first provide the proof for the 1-layer model. The values of $\mathcal{L}(\theta^\sharp)$, $\mathcal{L}(\theta^\sharp + \Delta\theta)$ are provided as

$$\mathcal{L}(\theta^\sharp) = \frac{1}{m} \times_{i \in \mathcal{S}} |\varepsilon_i|, \quad \mathcal{L}(\theta^\sharp + \Delta\theta) = \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| + \frac{1}{m} \times_{i \in \mathcal{S}^c} |\langle x_i, \Delta\theta \rangle - \varepsilon_i|. \quad (14)$$

Here $\mathcal{S} \in [m] : \{1, 2, \dots, m\}$ is the support of the noise vector $\varepsilon = [\varepsilon_1, \dots, \varepsilon_m]^\top$, and $\mathcal{S}^c = [m] - \mathcal{S}$. Hence, we have

$$\mathcal{L}(\theta^\sharp + \Delta\theta) - \mathcal{L}(\theta^\sharp) = \frac{1}{m} \times_{i \in \mathcal{S}} (|\langle x_i, \Delta\theta \rangle - \varepsilon_i| - |\varepsilon_i|) + \frac{1}{m} \times_{i \in \mathcal{S}^c} |\langle x_i, \Delta\theta \rangle|. \quad (15)$$

For simplicity, we denote $\mathbf{X} = [x_1, \dots, x_m]^\top \in \mathbb{R}^{m \times d}$ by the design matrix. To upper bound the infimum of $\mathcal{L}(\theta^\sharp + \Delta\theta) - \mathcal{L}(\theta^\sharp)$ over the ℓ_1 -norm ball $\{\Delta\theta : \|\Delta\theta\|_1 \leq \gamma\}$, we choose a specific $\overline{\Delta\theta}$ for each realization $\{(x_i, y_i)\}_{i=1}^m$ and show that $\overline{\Delta\theta}$ lies in the ℓ_1 -norm ball with high probability. To this goal, we pick $\overline{\Delta\theta}$ by solving the following linear equation

$$\mathbf{X}\Delta\theta = b, \quad (16)$$

where $b[i] = \text{Sign}(\varepsilon_i) (\xi \wedge |\varepsilon_i|)$ for $i \in \mathcal{S}$, and $b[i] = 0$ otherwise. Here $\xi > 0$ is a hyperparameter to be tuned later. Before showing that such $\overline{\Delta\theta}$ lies in the ℓ_1 -norm ball with high probability, we first further upper bound $\mathcal{L}(\theta^\sharp + \Delta\theta) - \mathcal{L}(\theta^\sharp)$ for this specific choice of $\overline{\Delta\theta}$. By our construction of $\overline{\Delta\theta}$, we first have

$$\begin{aligned} \mathcal{L}(\theta^\sharp + \overline{\Delta\theta}) - \mathcal{L}(\theta^\sharp) &= \frac{1}{m} \times_{i \in \mathcal{S}} |x_i, \overline{\Delta\theta} - \varepsilon_i| - |\varepsilon_i| + \frac{1}{m} \times_{i \in \mathcal{S}^c} |x_i, \overline{\Delta\theta}| \\ &= -\frac{1}{m} \times_{i \in \mathcal{S}} \xi \wedge |\varepsilon_i|. \end{aligned} \quad (17)$$

On the other hand, based on Assumption 1, each $\varepsilon_i, i \in \mathcal{S}$ has nonzero probability to be away from zero. Hence, we can define the event $\mathcal{E} := \{\text{There are at least } p_0 p m \text{ elements of } |\varepsilon| \text{ larger than } t_0\}$. Define $\mathcal{M}_0 = \{i : |\varepsilon_i| \geq t_0\}$. Using the tail bound of binomial distribution, we have $\mathbb{P}(\mathcal{E}) \geq \frac{1}{4}$. Conditioned on the event \mathcal{E} , we have $m_0 = |\mathcal{M}_0| \geq p_0 p m$ and

$$\mathcal{L}(\theta^\dagger + \overline{\Delta\theta}) - \mathcal{L}(\theta^\dagger) = -\frac{1}{m} \sum_{i \in \mathcal{S}} \xi \wedge |\varepsilon_i| \leq -p_0 p \cdot (\xi \wedge t_0). \quad (18)$$

Now it suffices to show that such $\overline{\Delta\theta}$ belongs to the ℓ_1 -norm ball with high probability. To this goal, we first choose a specific $\overline{\Delta\theta}$ satisfying $\mathbf{X} \overline{\Delta\theta} = b$. We divide $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$ where $\mathbf{X}_1 \in \mathbb{R}^{m \times m}$ and $\mathbf{X}_2 \in \mathbb{R}^{m \times (d-m)}$. Note that \mathbf{X}_1 is non-singular almost surely. Hence, we can pick $\overline{\Delta\theta} = [b^\top \mathbf{X}_1^{-1} (\mathbf{X}_1^\top \mathbf{X}_1)^{-1}, \mathbf{0}]^\top$. Next, it suffices to verify that with high probability, $\|\overline{\Delta\theta}\|_1 \leq \gamma$, which is equivalent to

$$(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top b \leq \gamma. \quad (19)$$

To this goal, note that with probability at least $1 - e^{-\Omega(m)}$

$$\begin{aligned} (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top b \leq (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top b & \\ \leq \frac{1}{m} \mathbf{X}_1^\top \mathbf{X}_1^{-1} \cdot \frac{1}{m} \mathbf{X}_1^\top \|b\| & \\ \stackrel{(a)}{\leq} \frac{1}{\sqrt{m}} \|b\| & \\ \leq \xi. & \end{aligned} \quad (20)$$

Here in (a) we used the fact that $\frac{1}{m} \mathbf{X}_1^\top \mathbf{X}_1^{-1} \leq 1$ and $\frac{1}{m} \mathbf{X}_1^\top \leq \frac{1}{m}$ with probability at least $1 - e^{-\Omega(m)}$. Hence, we only need to set $\xi \leq \gamma$, so that the chosen $\overline{\Delta\theta}$ belongs to the ℓ_1 -norm ball with probability at least $1 - e^{-\Omega(m)}$.

Combining the above steps and choosing $\xi \leq \gamma$, we know that with probability at least $\frac{1}{16}$, the below upper bound holds

$$\begin{aligned} \inf_{k \Delta} \inf_{k_1} \{\mathcal{L}(\theta^\dagger + \Delta\theta) - \mathcal{L}(\theta^\dagger)\} & \stackrel{\text{w.p. } 1}{\leq} \stackrel{(m)}{\leq} \mathcal{L}(\theta^\dagger + \overline{\Delta\theta}) - \mathcal{L}(\theta^\dagger) \\ & \stackrel{\text{w.p. } \frac{1}{4}}{\leq} -p_0 p \xi \\ & \leq -p_0 p \gamma. \end{aligned} \quad (21)$$

This completes the proof for the 1-layer model. For general N -layer models, we consider the true solution $\mathbf{w} \in \mathcal{W}$ with $w_1 = \theta^\dagger$, and $w_2 = \dots = w_N = \mathbf{1}$. Moreover, for $\Delta\mathbf{w}$ we choose $\Delta w_2 = \dots = \Delta w_N = \mathbf{0}$. It is easy to verify that

$$\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w} + \Delta\mathbf{w}) = \frac{1}{m} \sum_{i \in \mathcal{S}} (| \langle x_i, \Delta w_1 \rangle - \varepsilon_i | - |\varepsilon_i|) + \frac{1}{m} \sum_{i \in \mathcal{S}} | \langle x_i, \Delta w_1 \rangle |. \quad (22)$$

Therefore, an argument similar to the 1-layer model can be used to write

$$\inf_{\mathbf{w} \in \mathcal{W}} \inf_{\mathbf{w}^0: k\mathbf{w}} \inf_{\mathbf{w}^0 k_1} \{\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}^0)\} \leq \inf_{k \Delta \mathbf{w}_1 k_1} \{\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w} + \Delta\mathbf{w})\} \leq -p_0 p \gamma, \quad (23)$$

with probability of $\frac{1}{16}$.

Under-parameterized Regime

Given $\mathbf{w} \in \mathcal{W}$ and any $\Delta\mathbf{w}$, consider $\mathbf{w}^\theta = \mathbf{w} + \Delta\mathbf{w}$ and define

$$\Delta\theta = \sum_{i=1}^N (\theta^\dagger)^{\frac{N-i}{N}} \odot \sum_{j_1: \dots: j_i} \Delta w_{j_1} \odot \dots \odot \Delta w_{j_i}. \quad (24)$$

We have

$$\begin{aligned}\mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w} + \Delta\mathbf{w}) &= \frac{1}{m} \times_{i \in \mathcal{S}} (|\langle x_i, \Delta\theta \rangle - \varepsilon_i| - |\varepsilon_i|) + \frac{1}{m} \times_{i \in \mathcal{S}^c} |\langle x_i, \Delta\theta \rangle| \\ &\geq \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| - \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle|.\end{aligned}\quad (25)$$

Hence, it suffices to show that, with probability of $1 - e^{-\Omega(d)}$,

$$\inf_{\Delta \in \mathbb{R}^d} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| - \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| \geq 0. \quad (26)$$

Note that the above inequality is invariant with respect to scaling. Hence, it suffices to show that it holds for arbitrary $\Delta\theta \in S^{d-1}$ where $S^{d-1} := \{x \in \mathbb{R}^d : \|x\| = 1\}$ is the standard sphere. Hence, it suffices to show

$$\begin{aligned}\inf_{\Delta \in S^{d-1}} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| - \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| \\ \geq \inf_{\Delta \in S^{d-1}} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| - \sup_{\Delta \in S^{d-1}} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| \geq 0.\end{aligned}\quad (27)$$

For the first term, applying Lemma 6, we have that with probability at least $1 - e^{-\Omega(d)}$

$$\inf_{\Delta \in S^{d-1}} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| \geq \frac{2}{\pi}(1-p) - \frac{(1-p)d}{m}. \quad (28)$$

Similarly, for the second part, with probability of at least $1 - e^{-\Omega(d)}$, we have

$$\sup_{\Delta \in S^{d-1}} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| \leq \frac{2}{\pi}p + \frac{pd}{m}. \quad (29)$$

Combining both parts, we have that with probability at least $1 - e^{-\Omega(d)}$

$$\inf_{\Delta \in S^{d-1}} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| - \sup_{\Delta \in S^{d-1}} \frac{1}{m} \times_{i \in \mathcal{S}} |\langle x_i, \Delta\theta \rangle| \geq \frac{2}{\pi}(1-2p) - 2 \frac{d}{m} \geq 0. \quad (30)$$

The last inequality follows from the fact that $m \gg \frac{d}{(1-2p)^2}$. This completes the proof.

C.2 Proof of Theorem 4

Given any $\Delta\mathbf{w} = [\Delta w_1, \dots, \Delta w_N]^>$, the following equality holds for any point $\mathbf{w} = \mathbf{w}^? + \Delta\mathbf{w}$ where $\mathbf{w}^? = [\sqrt[?]{\theta^?}, \dots, \sqrt[?]{\theta^?}]^>$:

$$\begin{aligned}(\sqrt[?]{\theta^?} + \Delta w_1) \odot \dots \odot (\sqrt[?]{\theta^?} + \Delta w_N) - \theta^? &= \underbrace{(\theta^?)^{\frac{N-i}{N}} \odot \dots \odot \Delta w_{j_1} \odot \dots \odot \Delta w_{j_i}}_{:=\Delta_1} \\ &\quad + \underbrace{\Delta w_1 \odot \dots \odot \Delta w_N}_{:=\Delta_2}.\end{aligned}\quad (31)$$

Hence, we have

$$\mathcal{L}(\mathbf{w}^?) - \mathcal{L}(\mathbf{w}^? + \Delta\mathbf{w}) = \frac{1}{m} \times_{i \in \mathcal{S}} (|\langle x_i, \Delta\theta_1 + \Delta\theta_2 \rangle - \varepsilon_i| - |\varepsilon_i|) + \frac{1}{m} \times_{i \in \mathcal{S}^c} |\langle x_i, \Delta\theta_1 + \Delta\theta_2 \rangle|. \quad (32)$$

For simplicity, we denote Θ_1 as the set of $\Delta\theta_1$ defined in (31) with $\|\Delta\mathbf{w}\|_1 \leq \gamma$. Similarly, Θ_2 is the set of $\Delta\theta_2$ defined in (31) with $\|\Delta\mathbf{w}\|_1 \leq \gamma$.

Lower bound. To prove the lower bound, one can write

$$\begin{aligned} \mathcal{L}(\mathbf{w}^\?) - \mathcal{L}(\mathbf{w}^\? + \Delta\mathbf{w}) &\geq \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 + \Delta\theta_2 \rangle| - \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 + \Delta\theta_2 \rangle| \\ &\geq \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{1}{m} \sum_{i=1}^n |\langle x_i, \Delta\theta_2 \rangle| \end{aligned} \quad (33)$$

Hence, we have

$$\begin{aligned} \inf_{k\Delta\mathbf{w}^\?} \{\mathcal{L}(\mathbf{w}^\?) - \mathcal{L}(\mathbf{w}^\? + \Delta\mathbf{w})\} &\geq \inf_{\Delta \in \mathcal{S}_{1,2\Theta_1}} \left(\frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| \right) \\ &\quad - \sup_{\Delta \in \mathcal{S}_{2,2\Theta_2}} \frac{1}{m} \sum_{i=1}^n |\langle x_i, \Delta\theta_2 \rangle|. \end{aligned} \quad (34)$$

First, we bound the term $\inf_{\Delta \in \mathcal{S}_{1,2\Theta_1}} \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle|$. It is easy to see that the vector $\Delta\theta_1$ is k -sparse and has the same sparsity pattern as $\theta^\?$. Therefore, according to Lemma 6, there exist universal constants $C, c > 0$ such that the following hold

$$\mathbb{P}^{\otimes} \sup_{\Delta \in \mathcal{S}_{1,2\Theta_1}} \frac{1}{m^\theta \|\Delta\theta_1\|} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{r}{\pi} \geq C \frac{r}{m^\theta} + \delta \leq e^{-cm^\theta}, \quad (35)$$

and

$$\mathbb{P} \sup_{\Delta \in \mathcal{S}_{1,2\Theta_1}} \frac{1}{m^{\theta\theta} \|\Delta\theta_1\|} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{r}{\pi} \geq C \frac{r}{m^{\theta\theta}} + \delta \leq e^{-cm^{\theta\theta}}. \quad (36)$$

Here $m^\theta = (1-p)m$, and $m^{\theta\theta} = pm$. The inequality (35) implies that

$$\mathbb{P}^{\otimes} \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| \geq \|\Delta\theta_1\| \left(\frac{r}{\pi} (1-p) - C \frac{(1-p)k}{m} - \delta_1 \right), \forall \Delta\theta_1 \in \Theta_1 \geq 1 - e^{-\frac{cm}{1-p}}. \quad (37)$$

Similarly, the inequality (36) leads to

$$\mathbb{P} \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| \leq \|\Delta\theta_1\| \left(\frac{r}{\pi} p + C \frac{pk}{m} + \delta_2 \right), \forall \Delta\theta_1 \in \Theta_1 \geq 1 - e^{-\frac{cm}{p}}. \quad (38)$$

Upon setting $\delta_1 = \frac{(1-p)k}{m}$ and $\delta_2 = \frac{pk}{m}$, with probability of $1 - e^{-\Omega(k)}$, for all $\Delta\theta_1 \in \Theta_1$, we have

$$\frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| \geq \frac{r}{\pi} (1-2p) - C^\theta \frac{r}{m} \geq 0. \quad (39)$$

In the last inequality we used the assumption $m \gg \frac{k}{(1-2p)^2}$. The above argument implies

$$\inf_{\Delta \in \mathcal{S}_{1,2\Theta_1}} \left(\frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| - \frac{1}{m} \sum_{i \in \mathcal{S}} |\langle x_i, \Delta\theta_1 \rangle| \right) \geq 0, \quad (40)$$

with probability of at least $1 - e^{-\Omega(k)}$. Now we turn to bound the second part $\sup_{\Delta \in \mathcal{S}_{2,2\Theta_2}} \frac{1}{m} \sum_{i=1}^n |\langle x_i, \Delta\theta_2 \rangle|$. To this goal, we first apply Lemma 6, which leads to

$$\mathbb{P} \sup_{\Delta \in \mathcal{S}_{2,2\Theta_2}} \frac{1}{m \|\Delta\theta_2\|} \sum_{i=1}^n |\langle x_i, \Delta\theta_2 \rangle| - \frac{r}{\pi} \geq C \frac{r}{m} + \delta \leq e^{-cm^2}. \quad (41)$$

Therefore, upon setting $\delta = \frac{d}{m}$, with probability of $1 - e^{-\Omega(d)}$, we have

$$\begin{aligned}
\sup_{\Delta \in \mathcal{Z}_{2\Theta_2}} \frac{1}{m} \sum_{i=1}^n |\langle x_i, \Delta \theta_2 \rangle| &\leq \sup_{\Delta \in \mathcal{Z}_{2\Theta_2}} \|\Delta \theta_2\| \leq \frac{r}{\pi} + (C+1) \frac{r}{m} \\
&\leq \frac{r}{m} \sup_{k \in \mathcal{K}_1} \|\Delta w_1 \odot \dots \odot \Delta w_N\| \\
&= \frac{r}{m} \sup_{j \in \mathcal{J}} \sum_{i \in \mathcal{I}} \Delta w_i[j] \\
&= \frac{d}{\sqrt{m}} \gamma^N.
\end{aligned} \tag{42}$$

Here $\Delta w_i[j]$ is the j -th element of Δw_i . Therefore, we conclude that

$$\inf_{k \in \mathcal{K}_1} \{\mathcal{L}(\mathbf{w}^\dagger) - \mathcal{L}(\mathbf{w}^\dagger + \Delta \mathbf{w})\} \geq -\frac{d}{\sqrt{m}} \gamma^N, \tag{43}$$

with probability of $1 - e^{-\Omega(k)}$, thereby completing the proof of the lower bound.

Upper bound. To this goal, we first define a restricted set of perturbation vectors

$$\mathcal{V} := \{\mathbf{v} : \|\mathbf{v}\|_1 \leq \gamma, v_i[j] = 0, \forall i \in [N], j \in \text{supp}(\theta^\dagger)\}, \tag{44}$$

where $\text{supp}(\theta^\dagger)$ is the support of θ^\dagger . Based on this definition, we have

$$\begin{aligned}
\inf_{k \in \mathcal{K}_1} \{\mathcal{L}(\mathbf{w}^\dagger) - \mathcal{L}(\mathbf{w}^\dagger + \Delta \mathbf{w})\} &\leq \inf_{\Delta \mathbf{w} \in \mathcal{V}} \{\mathcal{L}(\mathbf{w}^\dagger) - \mathcal{L}(\mathbf{w}^\dagger + \Delta \mathbf{w})\} \\
&= \inf_{\Delta \mathbf{w} \in \mathcal{V}} \frac{1}{m} \sum_{i \in \mathcal{S}} (|\langle x_i, \Delta \theta_2 \rangle - \varepsilon_i| - |\varepsilon_i|) + \frac{1}{m} \sum_{i \in \mathcal{S}^c} |\langle x_i, \Delta \theta_2 \rangle|.
\end{aligned} \tag{45}$$

where $\Delta \theta_2 = \Delta w_1 \odot \dots \odot \Delta w_N$ is the same as before. Note that for any $\|\Delta \mathbf{w}\|_1 \leq \gamma$, we have $\|\Delta \theta_2\| \leq \sqrt{d} \gamma^N$. Moreover, this bound is attainable when $\Delta w_i \equiv [\pm \gamma, \dots, \pm \gamma]^T$. Hence, we have

$$\begin{aligned}
&\inf_{\Delta \mathbf{w} \in \mathcal{V}} \frac{1}{m} \sum_{i \in \mathcal{S}} (|\langle x_i, \Delta \theta_2 \rangle - \varepsilon_i| - |\varepsilon_i|) + \frac{1}{m} \sum_{i \in \mathcal{S}^c} |\langle x_i, \Delta \theta_2 \rangle| \\
&\leq \inf_{k \in \mathcal{K}_1} \frac{1}{m} \sum_{i \in \mathcal{S}} (|\langle x_i, \Delta \theta_2 \rangle - \varepsilon_i| - |\varepsilon_i|) + \frac{1}{m} \sum_{i \in \mathcal{S}^c} |\langle x_i, \Delta \theta_2 \rangle|.
\end{aligned} \tag{46}$$

An argument similar to the proof of Theorem 1 can be used to show that, with probability of at least $1/16$, we have

$$\inf_{k \in \mathcal{K}_1} \frac{1}{m} \sum_{i \in \mathcal{S}} (|\langle x_i, \Delta \theta_2 \rangle - \varepsilon_i| - |\varepsilon_i|) + \frac{1}{m} \sum_{i \in \mathcal{S}^c} |\langle x_i, \Delta \theta_2 \rangle| \leq \frac{p_0 p(d-k)}{m} \sqrt{d} \gamma^N. \tag{47}$$

Recalling that $k \ll d$, we have with probability of $1/16$

$$\inf_{k \in \mathcal{K}_1} \{\mathcal{L}(\mathbf{w}^\dagger) - \mathcal{L}(\mathbf{w}^\dagger + \Delta \mathbf{w})\} \geq -\sqrt{p_0 p} \frac{d}{\sqrt{m}} \gamma^N. \tag{48}$$

This completes the proof.

D Proofs of Convergence Analysis

D.1 Proof of Theorem 2

For simplicity of notation, we denote $u = w_1$ and $v = w_2$. Moreover, without loss of generality, we assume that the elements of θ° are arranged in descending order, i.e., $\theta_1^\circ \geq \dots \geq \theta_k^\circ > \theta_{k+1}^\circ = \dots = \theta_d^\circ = 0$, and the initial point satisfies $u_i, v_i = \Theta(\sqrt{\alpha})$, $\forall i \in [d]$. Moreover, for $v \in \mathbb{R}^d$, we define $v_{:i} = [v_1, \dots, v_i]^\top$ and $v_i = [v_i, \dots, v_d]^\top$. For short, we denote $v_{-i} = v_{i+1:}$. Finally, we define $\kappa = \theta_1^\circ / \theta_k^\circ$ as the condition number. Moreover, without loss of generality, we assume that $\theta_1^\circ \geq 1 \geq \theta_k^\circ$.

First, according to Proposition 1, the sub-differential of $\mathcal{L}(u, v)$ is uniformly concentrated around its population gradient. In particular, with probability at least $1 - Ce^{-cm^2}$, we have

$$u_{i;t+1} = u_{i;t} + \eta \frac{\theta_i^\circ - u_{i;t}v_{i;t}}{\|u_t \odot v_t - \theta^\circ\|} v_{i;t} + \eta \delta_i v_{i;t}, \text{ and } |\delta_i| \leq \delta, \forall i \in [d], \quad (49)$$

$$v_{i;t+1} = v_{i;t} + \eta \frac{\theta_i^\circ - u_{i;t}v_{i;t}}{\|u_t \odot v_t - \theta^\circ\|} u_{i;t} + \eta \delta_i u_{i;t}, \text{ and } |\delta_i| \leq \delta, \forall i \in [d]. \quad (50)$$

Hence, we have

$$u_{i;t+1}v_{i;t+1} = u_{i;t}v_{i;t} + \eta \frac{\theta_i^\circ - u_{i;t}v_{i;t}}{\|u_t \odot v_t - \theta^\circ\|} + \delta_i \left(u_{i;t}^2 + v_{i;t}^2 + \eta^2 \frac{\theta_i^\circ - u_{i;t}v_{i;t}}{\|u_t \odot v_t - \theta^\circ\|} + \delta_i \right) u_{i;t}v_{i;t}. \quad (51)$$

Moreover, we have

$$u_{i;t+1}^2 + v_{i;t+1}^2 = u_{i;t}^2 + v_{i;t}^2 \left(1 + \eta^2 \frac{(\theta_i^\circ - u_{i;t}v_{i;t})}{\|u_t \odot v_t - \theta^\circ\|} + \delta_i \right)^2 + 4\eta \frac{\theta_i^\circ - u_{i;t}v_{i;t}}{\|u_t \odot v_t - \theta^\circ\|} + \delta_i u_{i;t}v_{i;t}. \quad (52)$$

Signal Dynamics

We first study the behavior of the signal term $S_t = u_{:k;t}v_{:k;t}$ for the first k components of the model $u_t \odot v_t$. We divide the dynamics into $k + 1$ stages. In the first k stages, each component $u_{i;t}v_{i;t}$ converges to θ_i° sequentially. Once all the components are close to the ground truth, the distance between signal term and the ground truth $\|S_t - \theta_{:k}^\circ\|$ will further decrease to $\mathcal{O}(\sqrt{d^2 m \alpha^2} \vee \eta \theta_1^\circ)$.

Stage 1: In this stage, the first component $u_1 v_1$ grows to $\theta_1 - \delta \|\theta^\circ\|$ within $\Theta\left(\frac{k}{\eta} \log \frac{1}{\delta}\right)$ iterations. At the initial point, we have $u_{1;0}v_{1;0} = \Theta(\alpha)$. For iteration $t + 1$, according to (51), we have

$$\begin{aligned} u_{1;t+1}v_{1;t+1} &\geq u_{1;t}v_{1;t} + \eta \frac{\theta_1^\circ - u_{1;t}v_{1;t}}{\|u_t \odot v_t - \theta^\circ\|} + \delta_1 \left(u_{1;t}^2 + v_{1;t}^2 \right) \\ &\geq u_{1;t}v_{1;t} + 2\eta \frac{\theta_1^\circ - u_{1;t}v_{1;t}}{\|u_t \odot v_t - \theta^\circ\|} + \delta_1 u_{1;t}v_{1;t}. \end{aligned} \quad (53)$$

We further divide our analysis into two substages. In the first substage, we have $u_1 v_1 \leq \theta_1^\circ / 2$. Note that $\|u_t \odot v_t - \theta^\circ\| = \mathcal{O}(\|\theta^\circ\|)$ and $|\delta_1| \leq \delta \frac{1}{\eta}$. Hence, (53) can be further simplified as

$$u_{1;t+1}v_{1;t+1} \geq \left(1 + \frac{\eta \theta_1^\circ}{4 \|\theta^\circ\|} \right) u_{1;t}v_{1;t}. \quad (54)$$

Therefore, this stage ends within $\mathcal{O}\left(\frac{k}{\eta} \log \frac{1}{\delta}\right)$ iterations. In the second substage, we have $u_1 v_1 \geq \theta_1^\circ / 2$. Upon defining $x_t = \theta_1^\circ - u_{1;t}v_{1;t}$, one can write

$$\begin{aligned} x_{t+1} &\leq \left(1 - 2\eta \frac{u_{1;t}v_{1;t}}{\|u_t \odot v_t - \theta^\circ\|} \right) x_t - 2\eta \delta_1 u_t \odot v_t \\ &\leq \left(1 - \frac{\eta \theta_1^\circ}{\|\theta^\circ\|} \right) x_t + \eta \delta \theta_1^\circ. \end{aligned} \quad (55)$$

Hence, within additional $\mathcal{O}(\|\theta^\varnothing\|/(\eta\theta_1^\varnothing))$ iterations, we have $u_{1:T_1}v_{1:T_1} = \theta_1^\varnothing \pm \delta\|\theta^\varnothing\|$. Overall, within $T_1 = \mathcal{O}\left(\frac{k^\varnothing k}{\varnothing} \log \frac{1}{\delta}\right)$ iterations, we have $u_{1:T_1}v_{1:T_1} = \theta_1^\varnothing \pm \delta\|\theta^\varnothing\|$. Now, we turn to show the lower bound on T_1 by analyzing the trajectory of $u_{1:t}^2 + v_{1:t}^2$. Due to (52), when $u_{1:t}^2 + v_{1:t}^2 \leq \frac{\varnothing}{2}$, we have

$$u_{1;t+1}^2 + v_{1;t+1}^2 \leq u_{1:t}^2 + v_{1:t}^2 + 1 + 10\frac{\eta\theta_1^\varnothing}{\|\theta^\varnothing\|}. \quad (56)$$

Hence, at least $\Omega\left(\frac{k^\varnothing k}{\varnothing} \log \frac{1}{\delta}\right)$ iterations are needed for $u_{1:t}^2 + v_{1:t}^2$ to be larger than $\frac{\varnothing}{2}$. Since $u_{1:t}^2 + v_{1:t}^2 \geq u_{1:t}v_{1:t}$, we immediately obtain $T_1 = \Omega\left(\frac{k^\varnothing k}{\varnothing} \log \frac{1}{\delta}\right)$.

Stage 2 to Stage k : In the next $k-1$ stages, each component $u_i v_i$ will converge to $\theta_i^\varnothing \pm \delta\|\theta^\varnothing\|$ sequentially. To show this, we use an inductive argument. In each stage i , we assume that the first $i-1$ components have already converged close to $\theta_j^\varnothing, \forall j \in [i-1]$. Hence, we have $\|u_t \odot v_t - \theta^\varnothing\| = \Theta\left(\theta_{(i-1)}^\varnothing\right)$. Repeating the procedure in Stage 1, we can show that, at stage i , $T_i = \mathcal{O}\left(\frac{\|\theta_{(i-1)}^\varnothing\|}{\varnothing}\right) \log \frac{1}{\delta}$ iterations are needed for $u_i v_i$ to converge to $\theta_i^\varnothing \pm \delta\|\theta^\varnothing\|$. Overall, after $T = T_1 + \dots + T_k = \mathcal{O}\left(\frac{k^{3-2}}{\varnothing}\right) \log \frac{1}{\delta}$ iterations, we have

$$u_{:k:T}v_{:k:T} - \theta_{:k:T}^\varnothing \leq \sqrt{k}\delta\|\theta^\varnothing\|. \quad (57)$$

Stage $k+1$: In the final stage, the signal term will quickly decrease to $\mathcal{O}\left(\sqrt{d^2 m \alpha^{1-\Theta(\cdot)}} \vee \eta\theta_1^\varnothing\right)$ within $T_{k+1} = \mathcal{O}\left(\frac{k^\varnothing k}{\varnothing}\right)$ iterations. To show this, we write

$$\theta_{:k}^\varnothing - S_{t+1} = \theta_{:k}^\varnothing - S_t - \eta \frac{u_{:k;t}^2 + v_{:k;t}^2}{\|u_t \odot v_t - \theta^\varnothing\|} \odot (\theta_{:k}^\varnothing - S_t) - \eta \frac{u_t^2 + v_t^2 + \eta^2}{\|u_t \odot v_t - \theta^\varnothing\|} \frac{\theta_{:k}^\varnothing - S_t}{\|u_t \odot v_t - \theta^\varnothing\|} + \eta^2 \odot S_t. \quad (58)$$

Here we denote $\delta = [\delta_1, \dots, \delta_k]^\top$. Note that $\|u_t \odot v_t - \theta^\varnothing\| \leq \|\theta_{:k}^\varnothing - S_t\| + \|E_t\|$. Moreover, based on our assumption, we have $\|E_t\| \leq \sqrt{d}\alpha^{1-\Theta(\cdot)}$. Finally, the balanced property implies that $|u_{i:t} - v_{i:t}| = \mathcal{O}\left(\alpha^{1-\Theta(\cdot)}\right)$ (this will be proven later). Hence, we have

$$\begin{aligned} \|\theta_{:k}^\varnothing - S_{t+1}\| &\leq (\theta_{:k}^\varnothing - S_t) \odot \left(1 - \eta \frac{u_t^2 + v_t^2}{\|u_t \odot v_t - \theta^\varnothing\|} + \eta^2 \frac{(\theta_{:k}^\varnothing - S_t) \odot S_t}{\|u_t \odot v_t - \theta^\varnothing\|^2}\right) + 4\eta\delta\|\theta^\varnothing\| \\ &\leq \|\theta_{:k}^\varnothing - S_t\| \left(1 - \frac{\eta}{2} \frac{\theta_k^\varnothing}{\|\theta_{:k}^\varnothing - S_t\| + \|E_t\|}\right) + 4\eta\delta\|\theta^\varnothing\| \\ &\leq \|\theta_{:k}^\varnothing - S_t\| - 0.25\eta\theta_k^\varnothing + 4\eta\delta\|\theta^\varnothing\| \\ &\leq \|\theta_{:k}^\varnothing - S_t\| - 0.1\eta\theta_k^\varnothing. \end{aligned} \quad (59)$$

Here, we used the fact that $\delta \leq \frac{1}{2}$. On the other hand, we have

$$\begin{aligned} \|S_{t+1} - S_t\| &\leq \frac{\eta}{\|u_t \odot v_t - \theta^\varnothing\|} (u_t^2 + v_t^2 \odot (\theta_{:k}^\varnothing - S_t) + 4\eta\delta\|\theta^\varnothing\|) \\ &\stackrel{(a)}{\leq} 2\eta\theta_1^\varnothing + 4\eta\delta\|\theta^\varnothing\| \\ &\leq 3\eta\theta_1^\varnothing, \end{aligned} \quad (60)$$

where in (a) we used Lemma 7. The above inequality indicates that the signal propagation in each step is upper bounded by $\mathcal{O}(\eta\theta_1^\varnothing)$. Hence, we conclude that within $T_{k+1} = \mathcal{O}\left(\frac{k^\varnothing k}{\varnothing}\right)$ iterations, we have $\|\theta_{:k}^\varnothing - S_t\| \leq \sqrt{d^2 m \alpha^{1-\Theta(\cdot)}} \vee \eta\theta_1^\varnothing$. Since $\delta \leq \frac{1}{2}$, the total iteration complexity is upper bounded by $T^0 = T_1 + \dots + T_{k+1} = \mathcal{O}\left(\frac{k^{3-2}}{\varnothing}\right) \log \frac{1}{\delta}$.

Residual Dynamics

Now, we analyze the residual dynamics. Instead of analyzing the dynamics of $u_{i,t}v_{i,t}$, we analyze its surrogate $u_{i,t}^2 + v_{i,t}^2$. Based on (52), we can naturally bound it as follows

$$u_{i,t+1}^2 + v_{i,t+1}^2 \leq u_{i,t}^2 + v_{i,t}^2 + 6\eta\delta u_{i,t}v_{i,t} \leq (1 + 3\eta\delta) u_{i,t}^2 + v_{i,t}^2. \quad (61)$$

Therefore, during the training process, we can bound the residual term as

$$u_{i,t}^2 + v_{i,t}^2 \leq \alpha (1 + \eta\delta)^{\mathcal{O}\left(\frac{k^3-2}{k} \log\left(\frac{1}{\epsilon}\right)\right)} \leq \alpha^{1 - \mathcal{O}(k^3-2)}. \quad (62)$$

Hence, we have

$$\|E_t\| \leq \prod_{i=k+1}^d (u_{i,t}^2 + v_{i,t}^2)^{\frac{1}{2}} \leq \sqrt{d} \alpha^{1 - \mathcal{O}(k^3-2)}. \quad (63)$$

Therefore, we conclude that within $\bar{T} = \mathcal{O}\left(\frac{k^3-2}{k} \log\left(\frac{1}{\epsilon}\right)\right)$ iterations, we have

$$\|u_{\bar{T}} \odot v_{\bar{T}} - \theta^?\| \leq \|\theta_{:k}^? - S_{\bar{T}}\| + \|E_{\bar{T}}\| \leq \sqrt{d^2 m} \alpha^{1 - \mathcal{O}(k^3-2)} \vee \eta\theta_1^?. \quad (64)$$

Therefore, with probability at least $1 - e^{-\mathcal{C} \log^2(m) \log(d) \log(k^2/k^2)}$, we have

$$\|u_{\bar{T}} \odot v_{\bar{T}} - \theta^?\| \leq \sqrt{d^2 m} \alpha^{1 - \Theta\left(\frac{k^2}{m(1-p)^2}\right)} \vee \eta\theta_1^?. \quad (65)$$

Long Escape Time

Based on the above analysis, it can be seen that, after \bar{T} iterations, the residual term is the dominant term, which may cause the algorithm to diverge, as captured by Figure 1. We now show that the residual term will not diverge within $T^0 = \frac{m(1-p)^2}{k} \bar{T}$. To this goal, recall that for $\forall k+1 \leq i \leq d$, we have

$$u_{i,t+1}^2 + v_{i,t+1}^2 \leq (1 + 3\eta\delta) u_{i,t}^2 + v_{i,t}^2. \quad (66)$$

Hence, for $t \leq \frac{1}{6} \log\left(\frac{1}{\epsilon}\right)$, we have

$$u_{i,t}^2 + v_{i,t}^2 \leq \alpha (1 + 3\eta\delta)^{\frac{1}{6} \log\left(\frac{1}{\epsilon}\right)} \leq \alpha e^{0.5 \log\left(\frac{1}{\epsilon}\right)} = \sqrt{\alpha}. \quad (67)$$

The proof is completed by noticing that $m = \tilde{\Omega}\left(\frac{k}{(1-p)^2}\right)$.

Balanced Property

To prove the balanced property, we directly calculate the dynamic of the difference $u_{i,t} - v_{i,t}$. One can write

$$u_{i,t+1} - v_{i,t+1} = (u_{i,t} - v_{i,t}) \left(1 - \eta \frac{\theta_i^? - u_{i,t}v_{i,t}}{\|u_t \odot v_t - \theta^?\|}\right) - \eta\delta_i. \quad (68)$$

Since $\theta_i^? - u_{i,t}v_{i,t} \geq 0$, we have

$$|u_{i,t+1} - v_{i,t+1}| \leq |u_{i,t} - v_{i,t}| (1 + \eta\delta), \quad (69)$$

for $0 \leq i \leq k$. We conclude that

$$|u_{i,t} - v_{i,t}| \leq \sqrt{\alpha} (1 + \eta\delta)^t \leq \alpha^{0.5 - \Theta(k^3-2)}. \quad (70)$$

for $\forall t \leq \frac{k^3-2}{k} \log\left(\frac{1}{\epsilon}\right)$. On the other hand, for $i \geq k$, we can write

$$|u_{i,t} - v_{i,t}| \leq \frac{\mathcal{O}\left(\frac{1}{u_{i,t}^2 + v_{i,t}^2}\right)}{\alpha^{0.5 - \Theta(k^3-2)}}. \quad (71)$$

The proof is completed by noticing that $m = \tilde{\Omega}\left(\frac{k}{(1-p)^2}\right)$.

Convergence in Under-parameterized Regime

In this section, we study the under-parameterized regime, where we assume that $m = \tilde{\Omega}(\frac{d}{(1-\rho)^2})$. The analysis of the signal term is the same as the over-parameterized regime and hence omitted for brevity. Here we only analyze the residual dynamic. One can write

$$E_{t+1} = E_t \left[1 - \eta \frac{u_{k+1:t}^2 + v_{k+1:t}^2}{\|u_t \odot v_t - \theta^\tau\|^2} + \eta \delta_{k+1} \frac{u_{k+1:t}^2 + v_{k+1:t}^2}{\|u_t \odot v_t - \theta^\tau\|^2} + \eta^2 \frac{-E_t}{\|u_t \odot v_t - \theta^\tau\|^2} + \delta_{k+1}^2 E_t \right]. \quad (72)$$

When the residual term becomes the dominant term, i.e., $\|\theta_{:k}^\tau - S_t\| \leq \|E_t\|$, we have the simplified dynamic

$$\begin{aligned} \|E_{t+1}\| &\leq \|E_t\| \left[1 - 0.5\eta \frac{u_{k+1:t}^2 + v_{k+1:t}^2}{\|u_t \odot v_t - \theta^\tau\|^2} + \eta \delta_{k+1} \frac{u_{k+1:t}^2 + v_{k+1:t}^2}{\|u_t \odot v_t - \theta^\tau\|^2} + 2\eta^2 \|E_t\|^3 + \delta^2 \|E_t\| \right] \\ &\stackrel{(a)}{\leq} \|E_t\| \left[1 - \frac{\eta E_t}{\|E_t\|} + 4\eta \delta \|E_t\| \right] \\ &\leq \left(1 - \frac{\eta}{\sqrt{d}} \|E_t\| + 4\eta \delta \|E_t\| \right) \|E_t\| \\ &\stackrel{(b)}{\leq} \left(1 - 0.5 \frac{\eta}{\sqrt{d}} \|E_t\| \right) \|E_t\|. \end{aligned} \quad (73)$$

Here in (a) we used the balanced property, which results in $u_{k+1:t}^2 + v_{k+1:t}^2 \asymp 2E_t$. Moreover, (b) is implied by the fact that $m \gg \frac{d}{(1-\rho)^2}$, which in turn implies $\delta \ll \frac{1}{\rho^2 d}$. Hence, we have

$$\|u_{t+1} \odot v_{t+1} - \theta^\tau\| \leq \left(1 - \Omega \frac{\eta}{\sqrt{d}} \|u_t \odot v_t - \theta^\tau\| \right) \|u_t \odot v_t - \theta^\tau\|. \quad (74)$$

Then, for $t \geq \bar{T}$, we have

$$\|u_t \odot v_t - \theta^\tau\| \leq \sqrt{d^2 m \alpha} \left(1 - \Omega \frac{\eta}{\sqrt{d}} \|u_t \odot v_t - \theta^\tau\| \right)^{t - \bar{T}} \vee \eta \theta_1^\tau. \quad (75)$$

D.2 Proof of Theorem 3

The proof of N -layer model is similar to that of 2-layer model. First, we study the signal dynamics, showing that the first k components $w_{j,j}^{(t)}$ converge to θ_j^τ sequentially for $1 \leq j \leq k$. We also prove that the residual term remains small along the optimization trajectory. Based on the SubGM update rule, one can write

$$w_i^{(t+1)} = w_i^{(t)} + \eta \frac{\theta^\tau - w_j^{(t)}}{\theta^\tau - w_j^{(t)}} \prod_{j \neq i} w_j^{(t)} + \eta \prod_{j \neq i} w_j^{(t)}, \text{ and } \|\cdot\|_\tau \leq \delta, \forall i \in [N]. \quad (76)$$

Signal Dynamics

Stage 1: In this stage, we show that $w_{i,1}$ will converge to $\theta^?$ within $T_1 = \Theta \frac{k^?k}{N^?} \alpha \frac{N-2}{N}$ iterations. We first prove the upper bound. According to the update rule, we have

$$\begin{aligned} w_{i,1}^{(t+1)} &= w_{i,1}^{(t)} + \eta \frac{\theta_1^? - w_{j,1}^{(t)}}{\theta^? - w_j^{(t)}} + \delta_1 A @ w_{j,1}^{(t)} A + \text{higher order terms of } \eta \\ &\geq w_{i,1}^{(t)} + \eta \frac{\theta_1^? - w_{j,1}^{(t)}}{\theta^? - w_j^{(t)}} + \delta_1 A @ w_{j,1}^{(t)} A \\ &\geq w_{i,1}^{(t)} + N\eta \frac{\theta_1^? - w_{j,1}^{(t)}}{\theta^? - w_j^{(t)}} - \delta A w_{i,1}^{(t)} \frac{2(N-1)}{N}. \end{aligned} \quad (77)$$

Here we use the fact that $w_{i,1}^{(t)} \leq \theta_1^?$ and $\eta \cdot \frac{1}{N} \frac{N-2}{N}$ so that we can drop the higher order terms. For brevity, we only show how we can drop the 2-th order term of η . The proof of the higher order terms is similar. One can write

$$\begin{aligned} \eta^2 \sum_{i \neq j} \sum_{k \neq i} w_{k,1}^{(t)} w_{k,1}^{(t)} w_{k,1}^{(t)} A &\stackrel{(a)}{\leq} \eta^2 (\theta_1^?)^{\frac{N-2}{N}} \sum_{i \neq j} \sum_{k \neq i} w_{k,1}^{(t)} w_{k,1}^{(t)} A \\ &\stackrel{(b)}{\leq} (N-1) \eta (\theta_1^?)^{\frac{N-2}{N}} \eta \sum_{i=1}^{N-1} \sum_{j \neq i} w_{j,1}^{(t)} A \\ &\stackrel{(c)}{\leq} \eta \sum_{i=1}^{N-1} \sum_{j \neq i} w_{j,1}^{(t)} A. \end{aligned} \quad (78)$$

Here in (a) we use the balanced property and the fact that $w_{i,1}^{(t)} \leq \theta_1^?$. Moreover, in (b), we use the rearrangement Inequality. Finally in (c) we use the assumption that $\eta \cdot N^{-1} \frac{N-2}{N}$. For simplicity, we denote $x_t = w_{i,1}^{(t)}$. Note that $\theta^? - w_j^{(t)} \leq \|\theta^?\|$. Hence, the dynamic can be simplified as

$$x_{t+1} \geq x_t + \frac{N\eta}{\|\theta^?\|} (\theta_1^? - \delta \|\theta^?\| - x_t) x_t \frac{2(N-1)}{N}. \quad (79)$$

We next show that $x_T \geq \theta_1^? - 2\delta \|\theta^?\|$ within $T_1 = \Theta \frac{k^?k}{N^?} \alpha \frac{N-2}{N}$ iterations provided that $x_0 = \Theta(\alpha)$. To this goal, we divide our analysis into two substages.

- $x_t \leq \frac{\theta_1^?}{2}$: In this substage, we assume that $x_t \leq \frac{\theta_1^?}{2}$. Hence, we can further simplify the dynamic as

$$x_{t+1} \geq x_t + 0.5N\eta \frac{\theta_1^?}{\|\theta^?\|} x_t \frac{2(N-1)}{N}. \quad (80)$$

Without loss of generality, we assume that $x_0 = \alpha$. Now we divide the interval $[\alpha, 0.5\theta_1^?]$ into a series of sub-intervals $\{\mathcal{I}_k\}$, where $\mathcal{I}_k = [2^k\alpha, 2^{k+1}\alpha)$. In each \mathcal{I}_k , the dynamic can be further simplified as

$$x_{t+1} \geq (1 + 0.5N\eta \frac{\theta_1^?}{\|\theta^?\|} 2^k \alpha \frac{N-2}{N}) x_t. \quad (81)$$

Therefore, the number of iterations that x_t spends in each interval \mathcal{I}_k is $\mathcal{O} \frac{k^?k}{N^?} 2^k \alpha \frac{N-2}{N}$. Hence, the total number of iterations is upper bounded by

$$\mathcal{O} \sum_{k=0}^{P-1} \frac{k^?k}{N^?} 2^k \alpha \frac{N-2}{N} = \mathcal{O} \frac{k^?k}{N^?} \alpha \frac{N-2}{N}.$$

- $x_t \geq \frac{\theta_1^?}{2}$: In this substage, we define $y_t = \theta_1^? - \delta \|\theta^?\| - x_t$. Via a similar trick, we can show that within additional $\mathcal{O} \frac{k^?k}{N} (\theta_1^?)^{\frac{2N-2}{N}}$ iterations, we have $x_t \geq \theta_1^? - 2\delta \|\theta^?\|$. Overall, after $T_1 = \Theta \frac{k^?k}{N} \alpha^{\frac{N-2}{N}}$ iterations, we have $\theta_1^? - 2\delta \|\theta^?\| \leq \overset{\circ}{w}_{i;1}^{(T_1)} \leq \theta_1^?$.

Stages 2 to k : Similarly, for component $\overset{\circ}{w}_{j;l}^{(t)}$, it takes $\mathcal{O} \frac{\|\overset{\circ}{w}_{j;l}^{(t)}\|}{N} \alpha^{\frac{N-2}{N}}$ iterations to attain $\theta_j^? - 2\delta \|\theta^?\|$. Overall, Stages 2 to k take $\Theta \frac{k^{\frac{3}{2}}}{N} \alpha^{\frac{N-2}{N}}$ iterations to terminate.

Stage $k+1$: In this stage, we take $S_t = \overset{\circ}{w}_{j;k}^{(t)}$. Hence, we have

$$\begin{aligned} \|\theta_{:k}^? - S_{t+1}\| &\leq (\theta_{:k}^? - S_t) \overset{\circ}{\mathbb{B}} \mathbf{1} - \eta \frac{\overset{\circ}{\mathbb{P}}_{i=1}^N \overset{\circ}{\mathbb{Q}}_{j \neq i} w_{j;1}^{(t)} \overset{2}{\mathbb{A}} \overset{1}{\mathbb{C}}}{\theta^? - \overset{\circ}{w}_j^{(t)}} \overset{1}{\mathbb{A}} + 4N\eta\delta\sqrt{k}(\theta_1^?)^{\frac{2(N-1)}{N}} \\ &\leq \|\theta_{:k}^? - S_t\| \left(1 - N\eta \frac{(\theta_k^?)^{\frac{2(N-1)}{N}}}{\|\theta_{:k}^? - S_t\| + \|E_t\|}\right) + 4N\eta\delta\sqrt{k}(\theta_1^?)^{\frac{2(N-1)}{N}} \quad (82) \\ &\leq \|\theta_{:k}^? - S_t\| - 0.5N\eta(\theta_k^?)^{\frac{2(N-1)}{N}} + 4N\eta\delta\sqrt{k}(\theta_1^?)^{\frac{2(N-1)}{N}} \\ &\leq \|\theta_{:k}^? - S_t\| - 0.1N\eta(\theta_k^?)^{\frac{2(N-1)}{N}}. \end{aligned}$$

Here we used the fact that $\theta^? - \overset{\circ}{w}_j^{(t)} \leq \|\theta_{:k}^? - S_t\| + \|E_t\| \leq 2\|\theta_{:k}^? - S_t\|$, and the assumption that $\delta \leq \frac{1}{N} \frac{1}{2} \frac{2N-2}{N}$. On the other hand, we have

$$\begin{aligned} \|S_{t+1} - S_t\| &\leq \overset{\mathcal{X}}{\mathbb{N}} \eta \overset{\circ}{\mathbb{A}} \frac{\theta_{:k}^? - S_t}{\theta^? - \overset{\circ}{w}_j^{(t)}} + \overset{1}{i;k} \overset{\circ}{\mathbb{A}} \overset{\mathbb{Y}}{\mathbb{A}} \overset{1}{j \neq i} \overset{\circ}{w}_{j;k}^{(t)} \overset{\mathbb{A}}{\mathbb{A}} \quad (83) \\ &\leq 2N\eta\sqrt{k}(\theta_1^?)^{\frac{2(N-1)}{N}}. \end{aligned}$$

Hence, we conclude that within $\mathcal{O} \frac{\rho_k}{N} \frac{1}{2} \frac{2(N-1)}{N}$ iterations, we have $\|\theta_{:k}^? - S_t\| \leq \sqrt{d^2 m \alpha} \vee N\eta(\theta_1^?)^{\frac{2(N-1)}{N}}$. Overall, the total iteration complexity is bounded by $\mathcal{O} \frac{k^{\frac{3}{2}}}{N} \alpha^{\frac{N-2}{N}}$.

Residual Dynamics

Similar to the 2-layer model, here we study the surrogate of the residual term $\overset{\mathbb{P}}{\mathbb{N}}_{i=1} w_{i;l}^{(t)2}$ for $l \geq k+1$. To this goal, we first notice that

$$\begin{aligned} \overset{\mathcal{X}}{\mathbb{N}}_{i=1} w_{i;l}^{(t+1)2} &= \overset{\mathcal{X}}{\mathbb{N}}_{i=1} w_{i;l}^{(t)2} + 2N\eta \overset{\circ}{\mathbb{A}} \frac{-\overset{\circ}{w}_{j;l}^{(t)}}{\theta^? - \overset{\circ}{w}_j^{(t)}} + \delta_{i;l} \overset{\circ}{\mathbb{A}} \overset{\mathbb{Y}}{\mathbb{A}} w_{j;l}^{(t)} \\ &\quad + \eta^2 \overset{\mathcal{X}}{\mathbb{N}}_{i=1} \overset{\circ}{\mathbb{A}} \frac{-\overset{\circ}{w}_{j;l}^{(t)}}{\theta^? - \overset{\circ}{w}_j^{(t)}} + \delta_{i;l} \overset{\circ}{\mathbb{A}} \overset{\mathbb{Y}}{\mathbb{A}} w_{j;l}^{(t)} \overset{1}{j \neq i} \overset{\circ}{w}_{j;l}^{(t)} \overset{\mathbb{A}}{\mathbb{A}} \quad (84) \\ &\leq \overset{\mathcal{X}}{\mathbb{N}}_{i=1} w_{i;l}^{(t)2} + 4N\eta\delta \overset{\mathbb{Y}}{\mathbb{A}} w_{j;l}^{(t)} \\ &\leq \overset{\mathcal{X}}{\mathbb{N}}_{i=1} w_{i;l}^{(t)2} + 4N\eta\delta \overset{\circ}{\mathbb{B}} \frac{\overset{\mathbb{P}}{\mathbb{N}}_{i=1} w_{i;l}^{(t)2} \overset{1}{\mathbb{A}} \frac{N}{2}}{N} \overset{\mathbb{C}}{\mathbb{A}}. \end{aligned}$$

Hence, once we set $z_t = \sum_{i=1}^N w_{i;l}^{(t)2}$, we have the following simplified dynamic

$$z_{t+1} \leq z_t + 4N\eta\delta \frac{z_t}{N}^{\frac{N}{2}}, \quad (85)$$

with $z_0 = \Theta N\alpha^{\frac{2}{N}}$. We claim that within $\mathcal{O}\left(\frac{1}{N}\right)$ iterations, we still have $z_t = \Theta N\alpha^{\frac{2}{N}}$. To show this, we suppose without loss of generality that $z_0 = N\alpha^{\frac{2}{N}}$, and define T as the first time that $z_T \geq 2N\alpha^{\frac{2}{N}}$. For any $0 \leq t \leq T-1$, we have

$$z_{t+1} \leq z_t + 4N\eta\delta 2^{\frac{N}{2}}\alpha. \quad (86)$$

We conclude that

$$T \geq \frac{N\alpha^{\frac{2}{N}}}{4N\eta\delta 2^{\frac{N}{2}}\alpha} = \frac{1}{4\delta 2^{\frac{N}{2}}\eta} \frac{1}{\alpha^{\frac{N-2}{N}}} \& \frac{1}{N\eta} \alpha^{\frac{N-2}{N}}. \quad (87)$$

Therefore, via a basic inequality, we have

$$\Upsilon \quad w_{i;l}^{(t)} \leq \frac{\sum_{i=1}^N w_{i;l}^{(t)2} \frac{1}{N}}{\alpha}. \quad (88)$$

Combining the analysis of both signal and residual terms, we conclude that within $\Theta \frac{1}{N}\alpha^{\frac{N-2}{N}}$ iterations, we have

$$\Upsilon \quad w_{i;l}^{(t)} - \theta_l^? \leq \sqrt{d^2 m \alpha} \vee (\eta \theta_1^?)^{\frac{2(N-1)}{N}}. \quad (89)$$

Long Time Guarantee

Similar to the proof of the 2-layer model, one can show that the residual term becomes the dominant term in the generalization error, and it stays in the order of α within $\Omega \frac{1}{N}\alpha^{\frac{N-2}{N}}$ iterations. The details are omitted for brevity.

Balanced Property

To prove the balanced property, we first study the dynamic of $w_{i;l}^{(t)} - w_{j;l}^{(t)}$, $\forall l \in [k], i, j \in [N]$. To this goal, we have

$$w_{i;l}^{(t+1)} - w_{j;l}^{(t+1)} = w_{i;l}^{(t)} - w_{j;l}^{(t)} \left(1 - \eta \frac{\theta_l^? - w_{j;l}^{(t)}}{\theta_l^? - w_j^{(t)}}\right) + \delta_l \sum_{f \notin i,j} w_{f;l}^{(t)}, \quad (90)$$

which in turn implies

$$w_{i;l}^{(t+1)} - w_{j;l}^{(t+1)} \leq w_{i;l}^{(t)} - w_{j;l}^{(t)} \left(1 - \eta \frac{\theta_l^? - w_j^{(t)}}{\theta_l^? - w_j^{(t)}}\right) + \delta_l \sum_{f \notin i,j} w_f^{(t)}. \quad (91)$$

If $w_{i;l}^{(t)} \leq \theta_l^? - \delta \|\theta^?\|$, the above inequality can be simplified as

$$w_{i;l}^{(t)} - w_{j;l}^{(t)} \leq w_{i;l}^{(0)} - w_{j;l}^{(0)} \cdot \alpha^{\frac{1}{N}}, \forall i, j \in [N], l \in [k]. \quad (92)$$

Once $w_{i;l}^{(t)} \geq \theta_l^? - \delta \|\theta^?\|$, we immediately have $w_{i;l}^{(t)} = \frac{1}{N} \bar{\theta}_l^? \pm \mathcal{O}(\sqrt{N\alpha})$. Then, we show that $w_{i;l}^{(t)}$ will stay close to $\frac{1}{N} \bar{\theta}_l^?$. To this goal, we first observe that $w_{i;l}^{(t+1)} - w_{i;l}^{(t)} = w_{i;l}^{(t+1)} - w_{j;l}^{(t)} - w_{j;l}^{(t)}$, $\forall i, j \in [N]$, which indicates that $w_{i;l}^{(t)}$ increases or decreases simultaneously. Hence, we conclude that $w_{i;l}^{(t)} - w_{j;l}^{(t)} \leq \delta \frac{1}{N} \bar{\theta}_l^?$.

For the residual term, we can derive a tighter bound. First, we have

$$w_{i;l}^{(t+1)} = w_{i;l}^{(t)} + 2\eta \frac{-w_{j;l}^{(t)}}{\theta^2 - w_j^{(t)}} + \delta_l \frac{w_{j;l}^{(t)}}{w_j^{(t)}} + \eta^2 \frac{-w_{j;l}^{(t)}}{\theta^2 - w_j^{(t)}} + \delta_l \frac{w_{j;l}^{(t)}}{w_j^{(t)}}. \quad (93)$$

Since we have already shown that $w_{i;l}^{(t)} \leq \alpha$, we further have

$$w_{i;l}^{(t+1)} \leq w_{i;l}^{(t)} + 4\eta\delta\alpha. \quad (94)$$

Therefore, one can write $w_{i;l}^{(t)} \leq w_{i;l}^{(0)} + 4\eta\delta\alpha \frac{N-1}{N} \leq \alpha \frac{2}{N}$, which in turn implies $w_{i;l}^{(t)} - w_{j;l}^{(t)} \leq w_{i;l}^{(t)} + w_{j;l}^{(t)} \leq \alpha^{1-N}$.

Convergence in Under-parameterized Regime

Similar to the 2-layer model, we consider the dynamic of $E_t = w_{i;k+1}^{(t)}$, which is characterized as follows

$$\begin{aligned} \|E_{t+1}\| &\leq E_t + \sum_{i=1}^N \eta \frac{-E_t}{\theta^2 - w_j^{(t)}} + \delta_{k+1} \frac{w_{j;k+1}^{(t)}}{w_j^{(t)}} \\ &\leq E_t + N\eta \frac{-E_t}{\theta^2 - w_j^{(t)}} + \delta E_t^{\frac{2(N-1)}{N}} \\ &\leq E_t + N\eta \frac{-E_t}{\|E_t\|} + \delta E_t^{\frac{2(N-1)}{N}} \\ &\leq \|E_t\| - N\eta d^{\frac{N-1}{N}} \|E_t\|^{\frac{2N-2}{N}} + N\eta\delta E_t^{\frac{2(N-1)}{N}} \\ &\leq \|E_t\| - N\eta d^{\frac{N-1}{N}} \|E_t\|^{\frac{2N-2}{N}} + N\eta\delta \|E_t\|^{\frac{2N-2}{N}} \\ &\leq \|E_t\| - N\eta d^{\frac{N-1}{N}} \|E_t\|^{\frac{2N-2}{N}}. \end{aligned} \quad (95)$$

The last inequality comes from the fact that $\delta \leq d^{\frac{N-1}{N}}$ since we assume $m \geq \frac{d^{\frac{2N-2}{N}}}{(1-\rho)^2}$. Hence, we have

$$\|E_t\| \leq \frac{1}{N\eta d^{\frac{N-1}{N}} (N-1)N(t-\bar{T}) + 1/\|E_{\bar{T}}\|}. \quad (96)$$

Since the residual term is the dominant term in the generalization error, we have

$$w_i^{(t)} - \theta^2 \leq \frac{w_i^{(\bar{T})} - \theta^2}{N\eta d^{\frac{N-1}{N}} (N-1)N(t-\bar{T}) + 1}, \quad (97)$$

which completes the proof.

E Proof of Proposition 1

First, we provide an upper bound of the covering number for the (k, ϑ) -approximate sparse unit ball. We defer a preliminary discussion on covering number to Appendix H.

Lemma 4. Let $\mathcal{T}_{k,\#} := \{u \in \mathbb{R}^d : u \text{ is } (k, \vartheta)\text{-approximate sparse, } \|u\| \leq 1\}$. Then its covering number $N(\mathcal{T}_{k,\#}, \varepsilon, \|\cdot\|)$ is upper bounded by

$$N(\mathcal{T}_{k,\#}, \varepsilon, \|\cdot\|) \leq \frac{ed}{k}^k \left(1 + \frac{4}{\varepsilon}\right)^k, \quad (98)$$

provided that $\varepsilon \geq \vartheta$.

The next lemma will play a crucial role in proving Proposition 1.

Lemma 5. *Suppose $x \in \mathbb{R}^d$ is a standard Gaussian vector, i.e., $x_j \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$, and the noise ε satisfies Assumption 1, then we have*

$$\varphi(u) = \frac{\mathbb{E}[\text{Sign}(\langle x, u \rangle + \varepsilon) \langle x, v \rangle]}{\frac{u}{\|u\|}, v} = \frac{r}{\pi} (1-p) + \frac{r}{\pi} p \mathbb{E} e^{-2(2kuk^2)}.$$

The proof of this lemma can be found in Appendix G.1. Now, we are ready to prove Proposition 1. Our goal is to show that for arbitrary $u \in \mathcal{A}$, the following inequality holds

$$\frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_i - \varphi(u) \frac{u}{\|u\|} \leq \delta \quad (99)$$

with probability at least $1 - Ce^{-cm^2}$ provided that $m \geq \frac{k \log(d) \log(R) \log(\frac{1}{\delta})}{(p)^2}$. Here we define $\mathcal{A} := \{u : r \leq \|u\| \leq R, u \text{ is } (k, \vartheta)\text{-approximate sparse}\}$, where $r \geq \frac{1}{dm/k\vartheta \log(1/\vartheta)}$. Moreover, we define $\mathcal{B} := \{u : r \leq \|u\| \leq R, \|u\|_0 \leq k\}$ and $\mathcal{C} := \{(u, u^\theta) : u \in \mathcal{A}, v \in \mathcal{B}, \|u - u^\theta\| \leq \zeta\}$. Here \mathcal{B} is the ζ -net of \mathcal{B} with $\zeta \leq r$. Finally, we define $\mathcal{D} := \{\pm \mathbf{e}_j\}_{j \in [d]}$, where \mathbf{e}_j forms the standard basis of \mathbb{R}^d . Based on these definitions, we have

$$\begin{aligned} & \sup_{u \in \mathcal{A}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_i - \varphi(u) \frac{u}{\|u\|} \\ &= \sup_{u \in \mathcal{A}, v \in \mathcal{D}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) \langle x_i, v \rangle - \frac{\varphi(u)}{\|u\|} \langle u, v \rangle \\ &= \sup_{v \in \mathcal{D}} \sup_{u \in \mathcal{A}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) \langle x_i, v \rangle - \frac{\varphi(u)}{\|u\|} \langle u, v \rangle. \end{aligned} \quad (100)$$

We then show that for each element $y \in \mathcal{D}$, $\sup_{u \in \mathcal{A}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) \langle x_i, y \rangle - \varphi(u) \frac{\langle u, y \rangle}{\|u\|}$, $j \in [d]$ is $\mathcal{O}(\frac{1}{m})$ -sub-Gaussian random variable. To see this, note that

$$\begin{aligned} \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{ij} - \varphi(u) \frac{u_j}{\|u\|} &\leq \|\text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{ij}\|_2 + \frac{r}{\pi} \\ &\leq \|x_{ij}\|_2 + \frac{r}{\pi} = \mathcal{O}(1). \end{aligned} \quad (101)$$

Here we use the property of sub-Gaussian norm. This implies that $\frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{ij} - \varphi(u) \frac{u_j}{\|u\|}$ is $\mathcal{O}(\frac{1}{m})$ -sub-Gaussian random variable, since it is the sample average of $\text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{ij} - \varphi(u) \frac{u_j}{\|u\|}$.

Hence, via maximal inequality, we have that for $\forall t > 0$,

$$\begin{aligned} & \mathbb{P} \left(\sup_{u \in \mathcal{A}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_i - \varphi(u) \frac{u}{\|u\|} \geq t \right) \\ & \geq \sup_{y \in \mathcal{D}} \mathbb{E} \left(\sup_{u \in \mathcal{A}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) \langle x_i, y \rangle - \varphi(u) \frac{\langle u, y \rangle}{\|u\|} \right) + t \\ & \leq 2de^{-cm^2}. \end{aligned} \quad (102)$$

Hence, it suffices to study $\mathbb{E} \left(\sup_{u \in \mathcal{A}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|} \right)$. To this goal, we decompose it into two terms via triangle inequality.

$$\mathbb{E} \left(\sup_{u \in \mathcal{A}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|} \right) \leq (A) + (B), \quad (103)$$

where

$$(A) := \mathbb{E} \sup_{u \in \mathcal{B}} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|}, \quad (104)$$

and

$$(B) := \mathbb{E} \sup_{(u, u^\theta) \in \mathcal{Z}} \frac{1}{m} \sum_{i=1}^m (\text{Sign}(\langle x_i, u \rangle + \varepsilon_i) - \text{Sign}(\langle x_i, u^\theta \rangle + \varepsilon_i)) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|} + \varphi(u^\theta) \frac{u_1^\theta}{\|u^\theta\|}. \quad (105)$$

We first control (A). To this goal, we apply the union bound. Note that $\frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|}$ is $\mathcal{O}(\frac{1}{m})$ -sub-Gaussian and $|\mathcal{B}| \leq R^{Ck \log(d)}$. We then have

$$(A) \leq \frac{\sqrt{Ck \log(d) \log \frac{R}{m}}}{m}. \quad (106)$$

Now we control (B). Via triangle inequality, we first obtain

$$(B) \leq \mathbb{E} \underbrace{\sup_{(u, u^\theta) \in \mathcal{Z}} \frac{1}{m} \sum_{i=1}^m (\text{Sign}(\langle x_i, u \rangle + \varepsilon_i) - \text{Sign}(\langle x_i, u^\theta \rangle + \varepsilon_i)) x_{i,1}}_{(B_1)} + \underbrace{\sup_{(u, u^\theta) \in \mathcal{Z}} -\varphi(u) \frac{u_1}{\|u\|} + \varphi(u^\theta) \frac{u_1^\theta}{\|u^\theta\|}}_{(B_2)}. \quad (107)$$

For the first part, applying Hölder's inequality leads to

$$\begin{aligned} (B_1) &\leq \mathbb{E} \sup_{(u, u^\theta) \in \mathcal{Z}} \frac{1}{m} \sum_{i=1}^m |\text{Sign}(\langle x_i, u \rangle + \varepsilon_i) - \text{Sign}(\langle x_i, u^\theta \rangle + \varepsilon_i)| \max_{1 \leq i \leq m} |x_{i,1}| \\ &\leq \mathbb{E} \sup_{(u, u^\theta) \in \mathcal{Z}} \frac{1}{m} \sum_{i=1}^m \mathbb{1}(|\langle x_i, u - u^\theta \rangle| \geq |\langle x_i, u \rangle + \varepsilon_i|) \max_{1 \leq i \leq m} |x_{i,1}| \\ &\leq \mathbb{E} \sup_{k \Delta u^k} \frac{1}{m} \sum_{i=1}^m \mathbb{1}(|\langle x_i, \Delta u \rangle| \geq t) \max_{1 \leq i \leq m} |x_{i,1}| \\ &\quad + \mathbb{E} \sup_{u \in \mathcal{B}} \frac{1}{m} \sum_{i=1}^m \mathbb{1}(|\langle x_i, u \rangle + \varepsilon_i| \leq t) \max_{1 \leq i \leq m} |x_{i,1}|, \end{aligned} \quad (108)$$

where $t > 0$ is a constant to be determined later. Here, we used the fact that $\mathbb{1}(|\langle x_i, u - u^\theta \rangle| \geq |\langle x_i, u \rangle + \varepsilon_i|) \leq \mathbb{1}(|\langle x_i, \Delta u \rangle| \geq t) + \mathbb{1}(|\langle x_i, u \rangle + \varepsilon_i| \leq t)$ in the last inequality. We first bound (B₃)

$$\begin{aligned} (B_3) &\leq \mathbb{E} \frac{1}{m} \sum_{i=1}^m \mathbb{1}(\zeta \|x_i\| \geq t) \max_{1 \leq i \leq m} |x_{i,1}| \\ &\leq \mathbb{E} [\mathbb{1}(\zeta \|x_i\| \geq t)] \mathbb{E} \max_{j \neq i} |x_{j,1}| + \mathbb{E} [\mathbb{1}(\zeta \|x_i\| \geq t) |x_{i,1}|] \\ &\leq e^{-\frac{t^2}{2}} \frac{1}{\log(m)} + \mathbb{E} [\mathbb{1}(\zeta \|x_i\| \geq t) |x_{i,1}|], \end{aligned} \quad (109)$$

provided that $t \geq \sqrt{d}$. Applying Cauchy-Schwarz inequality, we have

$$\mathbb{E} [\mathbb{1}(\zeta \|x_i\| \geq t) |x_{i,1}|] \leq \sqrt{\frac{\mathbb{P}(\zeta \|x_i\| \geq t)}{\mathbb{P}(\zeta \|x_i\| \geq t)}} \sqrt{\mathbb{E} x_{i,1}^2} \leq e^{-\frac{t^2}{2}}. \quad (110)$$

Hence, we conclude that (B₃) $\leq e^{-c \frac{t^2}{2} \rho \frac{1}{\log(m)}}$. Next we control (B₄). Note that $\max_i |x_{i,1}|$ is $\mathcal{O}(\log(m))$ -sub-Gaussian. Via union bound, we have

$$(B_4) \leq \sup_{u \in B} \mathbb{E} \mathbb{1}(|\langle x_i, u \rangle + \varepsilon_i| \leq t) \max_{i=1}^m |x_{i,1}| + C \frac{\sqrt{k \log(m) \log(d) \log(\frac{R}{\#})}}{m}. \quad (111)$$

For the first part, applying the similar decomposition method, we have

$$\begin{aligned} \mathbb{E} \mathbb{1}(|\langle x_i, u \rangle + \varepsilon_i| \leq t) \max_{i=1}^m |x_{i,1}| &\leq \mathbb{E} \mathbb{1}(|\langle x_i, u \rangle + \varepsilon_i| \leq t) \max_{j \neq i} |x_{j,1}| \\ &+ \mathbb{E} [\mathbb{1}(|\langle x_i, u \rangle + \varepsilon_i| \leq t) |x_{i,1}|] \\ &\leq \rho \frac{t}{\log(m)} + \frac{1}{r}. \end{aligned} \quad (112)$$

Hence, we conclude that (B₄) $\leq \rho \frac{t}{\log(m)} + \frac{1}{r} + \frac{C k \log(m) \log(d) \log(\frac{R}{\#})}{m}$. For (B₂), we first have

$$\begin{aligned} -\varphi(u) \frac{u_1}{\|u\|} + \varphi(u^\theta) \frac{u_1^\theta}{\|u^\theta\|} &= |\varphi(u^\theta) - \varphi(u)| \frac{|u_1|}{\|u\|} + \varphi(u^\theta) \frac{u_1^\theta}{\|u^\theta\|} - \frac{u_1}{\|u\|} \\ &\leq |\varphi(u^\theta) - \varphi(u)| + \zeta. \end{aligned} \quad (113)$$

For the first part, we use Mean Value Theorem to write

$$|\varphi(u^\theta) - \varphi(u)| \leq \|\nabla \varphi(v)\| \|u^\theta - u\| \leq \|\nabla \varphi(v)\| \zeta, \quad (114)$$

where v is a point between u and u^θ . Note that $\nabla \varphi(v) = \frac{1}{2} \rho \mathbb{E} \frac{v^2}{\|v\|^4} e^{-\frac{v^2}{2\kappa v \kappa^2}}$. Hence, we have

$$\sup_{\|v\| \leq r} \|\nabla \varphi(v)\| \leq \sup_{\|v\| \leq r} \mathbb{E} \frac{e^2}{\|v\|^3} e^{-\frac{v^2}{2\kappa v \kappa^2}} \leq \frac{1}{r} \sup_{\|v\| \leq r} \mathbb{E} \frac{e^2}{\|v\|^2} e^{-\frac{v^2}{2\kappa v \kappa^2}} \leq \frac{1}{r}. \quad (115)$$

Overall, we have (B₂) $\leq \frac{1}{r}$, which results in

$$\begin{aligned} \mathbb{E} \sup_{u \in A} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|} &\leq \frac{\zeta}{r} + e^{-c \frac{t^2}{2} \rho \frac{1}{\log(m)}} + \rho \frac{t}{\log(m)} + \frac{C k \log(m) \log(d) \log(\frac{R}{\#})}{m}. \end{aligned} \quad (116)$$

Hence, once we set $\zeta \asymp \vartheta$, and $t \asymp \sqrt{d} \vartheta \log(m)$, together with the assumption that $r \asymp \frac{dm}{k} \vartheta \log \frac{1}{\#}$, we conclude that

$$\mathbb{E} \sup_{u \in A} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|} \leq \frac{C k \log^2(m) \log(d) \log(\frac{R}{\#})}{m}. \quad (117)$$

This leads to

$$\begin{aligned} \mathbb{P} \left(\sup_{u \in A} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \varphi(u) \frac{u_1}{\|u\|} \geq C \frac{k \log^2(m) \log(d) \log(\frac{R}{\#})}{m} + \delta \right) &\leq 2de^{-cm^2}. \end{aligned} \quad (118)$$

Therefore, the following inequality holds, provided that $m \asymp \frac{k \log^2(m) \log(d) \log(\frac{R}{\#})}{(1-\rho)^2}$

$$\mathbb{P} \left(\sup_{u \in A} \frac{1}{\varphi(u)} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i,1} - \frac{u}{\|u\|} \geq \delta \right) \leq e^{-cm^2}. \quad (119)$$

Now we turn to the case $m \ll \frac{d}{(1-p)^2}$. Following the same technique, it suffices to bound $\mathbb{E} \sup_{u \in \mathbb{R}^d} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i:1} - \varphi(u) \frac{u_1}{\|u\|}$. To this goal, we first notice that

$$\begin{aligned} & \mathbb{E} \sup_{u \in \mathbb{R}^d} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_{i:1} - \varphi(u) \frac{u_1}{\|u\|} \\ &= \mathbb{E} \sup_{\|u\|=1} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \lambda \varepsilon_i) x_{i:1} - \varphi(u) u_1. \end{aligned} \quad (120)$$

Similarly, applying one-step discretization, we have

$$\begin{aligned} (A) &\leq \mathbb{E} \sup_{u \in \mathcal{S}^d} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \lambda \varepsilon_i) x_{i:1} - \phi(\lambda) u_1 \\ &+ \mathbb{E} \sup_{\substack{u, u^0 \in \mathcal{S}^d \\ \|u - u^0\| \leq \frac{1}{m}}} \frac{1}{m} \prod_{i=1}^m (\text{Sign}(\langle x_i, u \rangle + \lambda \varepsilon_i) - \text{Sign}(\langle x_i, u^0 \rangle + \lambda \varepsilon_i)) x_{i:1} + \phi(\lambda) (u_1^0 - u_1). \end{aligned} \quad (121)$$

Here $\phi(\lambda) = \frac{1}{2}(1-p) + \frac{1}{2}p \mathbb{E} e^{-\lambda^2} = \frac{1}{2}(1-p) + \frac{1}{2}p e^{-\frac{1}{2}}$ is the same as before. We first control (B). To this goal, we show that $\sup_{u \in \mathbb{R}^d} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \lambda \varepsilon_i) x_{i:1} - \phi(\lambda) u_1$ is $\mathcal{O}(1/m)$ -sub-Gaussian. We prove it via checking the sub-Gaussian norm

$$\sup_{u \in \mathbb{R}^d} \text{Sign}(\langle x_i, u \rangle + \lambda \varepsilon_i) x_{i:1} - \phi(\lambda) u_1 \leq \|x_{i:1}\|_2 + \frac{1}{\pi} = \mathcal{O}(1). \quad (122)$$

Hence, via maximum inequality, we have

$$(B) \leq \mathbb{E} \sup_{u \in \mathbb{R}^d} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \lambda \varepsilon_i) x_{i:1} - \phi(\lambda) u_1 + \mathcal{O}\left(\frac{d \log \frac{1}{m}}{m}\right). \quad (123)$$

To control (D), we further decompose it into two parts,

$$\begin{aligned} (D) &\leq \mathbb{E} \sup_{\nu \in [0,1]} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\nu \langle x_i, u \rangle + \varepsilon_i) x_{i:1} - \phi\left(\frac{1}{\nu}\right) u_1 \\ &+ \mathbb{E} \sup_{\nu \in [0,1]} \frac{1}{m} \prod_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \lambda \varepsilon_i) x_{i:1} - \phi(\lambda) u_1. \end{aligned} \quad (124)$$

To control (D₁) and (D₂) we use arguments based on bracketing maximal inequality. We defer a preliminary discussion on bracketing maximal inequality to Appendix H. We first control (D₁). Let \mathcal{T} be defined as the ξ -net of the interval [0, 1]. We show that for any $\nu, \nu^0 \in [0, 1]$ such that $|\nu - \nu^0| \leq \xi$, we can control $\|(\text{Sign}(\nu \langle x_i, u \rangle + \varepsilon_i) - \text{Sign}(\nu^0 \langle x_i, u \rangle + \varepsilon_i)) x_{i:1}\|_{L_2(\mathcal{P})}$. To this goal, we first have

$$\begin{aligned} & \mathbb{E} \left[(\text{Sign}(\nu \langle x_i, u \rangle + \varepsilon_i) - \text{Sign}(\nu^0 \langle x_i, u \rangle + \varepsilon_i))^2 x_{i:1}^2 \right] \\ & \leq \mathbb{E} [|\text{Sign}(\nu \langle x_i, u \rangle + \varepsilon_i) - \text{Sign}(\nu^0 \langle x_i, u \rangle + \varepsilon_i)|] \\ & \leq \mathbb{E} [\mathbb{1}(|\nu - \nu^0| \langle x_i, u \rangle + \varepsilon_i| \geq t) + \mathbb{1}(|\nu \langle x_i, u \rangle + \varepsilon_i| \leq t)] \\ & \leq e^{-\frac{t^2}{2}} + t. \end{aligned} \quad (125)$$

Upon picking $t \asymp \xi \log \frac{1}{\xi}$, we have

$$\|(\text{Sign}(\nu \langle x_i, u \rangle + \varepsilon_i) - \text{Sign}(\nu^0 \langle x_i, u \rangle + \varepsilon_i)) x_{i:1}\|_{L_2(\mathbb{P})} \leq \frac{1}{\xi \log \frac{1}{\xi}}. \quad (126)$$

Therefore, the bracketing number is bounded by $N_{[]}(\varepsilon \|F\|, \mathcal{F}, \|\cdot\|) \leq C \varepsilon^{-\frac{1}{\pi}}$, which in turn leads to an upper bound on the bracketing entropy $J_{[]}^*(1, \mathcal{F}, L_2(\mathbb{P})) \leq 1$. Applying Theorem 6 leads to

$$(D_1) \leq \frac{1}{m}. \quad (127)$$

Similarly, we can show that $(D_2) \leq \frac{1}{m}$. Therefore, we conclude that

$$(B) \leq \frac{d \log \frac{1}{\varepsilon}}{m}. \quad (128)$$

For (C), we can use the similar technique in the overparameterized setting ($m \ll d$), which leads to

$$(C) \leq \frac{1}{\log(m)\varepsilon} + \frac{d \log(m)}{m}. \quad (129)$$

Therefore, once we set $\varepsilon \asymp \frac{d}{m}$, we immediately obtain

$$(A) \leq \frac{d \log(m)}{m}. \quad (130)$$

Combining the derived bounds results in

$$\mathbb{P} \sup_{u \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_i - \varphi(u) \frac{u}{\|u\|} \geq C \frac{d \log(m)}{m} + \delta \leq e^{c_1 \log(d)} c_2 m^{-2}. \quad (131)$$

Assuming $m \geq \frac{d \log(m)}{(1-\rho)^2}$, the above bound reduces to

$$\mathbb{P} \sup_{u \in \mathbb{R}^d} \frac{1}{\varphi(u)} \frac{1}{m} \sum_{i=1}^m \text{Sign}(\langle x_i, u \rangle + \varepsilon_i) x_i - \frac{u}{\|u\|} \geq \delta \leq e^{-cm^2}. \quad (132)$$

F Auxiliary Lemmas

Lemma 6. *Suppose x_1, \dots, x_m are i.i.d. standard Gaussian vectors with dimension d . Then, for arbitrary $\delta > 0$ we have*

$$\mathbb{P} \sup_{\|u\|=1} \frac{1}{m} \sum_{i=1}^m |\langle x_i, u \rangle| - \frac{2}{\pi} \geq C \frac{d}{m} + \delta \leq e^{-cm^2}. \quad (133)$$

Here C, c are universal constants.

Proof. This lemma directly follows from the standard expectation and high probability bounds for sub-Gaussian process. See e.g., [29, Lemma 4] for a simple proof. \square

Lemma 7. *For two arbitrary vectors $a, b \in \mathbb{R}^n$, we have*

$$\|a \odot b\| \leq \|a\|_1 \|b\|. \quad (134)$$

G Deferred Proofs

G.1 Proof of Lemma 5

Proof. To prove this lemma, it suffices to show that, for any $u, v \in \mathbb{R}^d$, we have

$$\mathbb{E} [\text{Sign}(\varepsilon + \langle x, u \rangle) \langle x, v \rangle] = \frac{1}{\pi} \mathbb{E} \left[e^{-2ku^2} \frac{u}{\|u\|}, v \right]. \quad (135)$$

Without loss of generality, we assume that $\|u\| = \|v\| = 1$. Let us denote $w := \langle x, u \rangle$, $z := \langle x, v \rangle$, $\rho := \text{Cov}(w, z) = \langle u, v \rangle$. Then

$$\begin{aligned} \mathbb{E} [\text{Sign}(\varepsilon + \langle x, u \rangle) \langle x, v \rangle] &= \mathbb{E} [\text{Sign}(\varepsilon + w) z] \\ &\stackrel{(a)}{=} \rho \mathbb{E} [\text{Sign}(w + \varepsilon) w] \\ &= \rho \mathbb{E} \left[\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt - \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \right] \\ &= \rho \mathbb{E} \left[\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt + \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \right] \\ &= 2\rho \mathbb{E} \left[\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \right] \\ &= \frac{1}{\pi} \langle u, v \rangle \mathbb{E} \left[\int_{-\infty}^{\infty} e^{-t^2/2} dt \right] \\ &= \frac{1}{\pi} \langle u, v \rangle \mathbb{E} \left[e^{-2ku^2} \right]. \end{aligned} \quad (136)$$

Here in (a) we use the fact that $z|w, \varepsilon \sim \mathcal{N}(\rho w, 1 - \rho^2)$ since ε is independent of w, z . Hence, we have

$$\mathbb{E} [\text{Sign}(\varepsilon + \langle x, u \rangle) \langle x, v \rangle] = \frac{1}{\pi} \mathbb{E} \left[e^{-2ku^2} \frac{u}{\|u\|}, v \right] \quad (137)$$

for any $u, v \in \mathbb{R}^d$. On the other hand, it is easy to verify that $\mathbb{E} [\text{Sign}(\langle x, u \rangle) \langle x, v \rangle] = \frac{1}{\pi} \frac{\langle u, v \rangle}{\|u\| \|v\|}$. The proof is completed by noting that the corruption probability is p . \square

H Preliminaries on the Uniform Concentration Bounds

In this section, we provide the preliminary probability tools for proving Proposition 1.

Definition 2 (Sub-Gaussian random variable). *We say a random variable $X \in \mathbb{R}$ with expectation $\mathbb{E}[X] = \mu$ is σ^2 -sub-Gaussian if for all $\lambda \in \mathbb{R}$, we have $\mathbb{E} e^{\lambda(X - \mu)} \leq e^{\frac{\lambda^2 \sigma^2}{2}}$. Moreover, the sub-Gaussian norm of X is defined as $\|X\|_2 := \sup_{\lambda \in \mathbb{R}} \lambda^{-1} \mathbb{E} [e^{\lambda(X - \mu)}]^{1/p}$.*

According to [39], the following statements are equivalent:

- X is σ^2 -sub-Gaussian.
- (Tail bound) For any $t > 0$, we have $\mathbb{P}(|X - \mu| \geq t) \leq 2e^{-\frac{t^2}{2\sigma^2}}$.
- (Moment bound) We have $\|X\|_2 \leq \sigma$.

Next, we provide the definitions of the sub-Gaussian process, ε -net, and covering number.

Definition 3 (Sub-Gaussian process). *A zero mean stochastic process $\{\mathcal{X}, \theta \in \mathbb{T}\}$ is a σ^2 -sub-Gaussian process with respect to a metric d on a set \mathbb{T} , if for every $\theta, \theta^0 \in \mathbb{T}$, the random variable $\mathcal{X} - \mathcal{X}^0$ is $(\sigma d(\theta, \theta^0))^2$ -sub-Gaussian.*

Definition 4 (ε -net and covering number). *A set \mathcal{N} is called an ε -net for (\mathbb{T}, d) if for every $t \in \mathbb{T}$, there exists $\pi(t) \in \mathcal{N}$ such that $d(t, \pi(t)) \leq \varepsilon$. The covering number $N(\mathbb{T}, d, \varepsilon)$ is defined as the smallest cardinality of an ε -net for (\mathbb{T}, d) :*

$$N(\mathbb{T}, d, \varepsilon) := \inf\{|\mathcal{N}| : \mathcal{N} \text{ is an } \varepsilon\text{-net for } (\mathbb{T}, d)\}.$$

Definition 5 (Bracketing number, Definition 2.1.6 in [37]). *Given two functions l and u , the bracket $[l, u]$ is the set of all functions f with $l \leq f \leq u$. An ε -bracket is a bracket $[l, u]$ with $\|u - l\| < \varepsilon$. The bracketing number $N_{[\cdot]}(\varepsilon, \mathcal{F}, \|\cdot\|)$ is the minimum number of ε -brackets needed to cover \mathcal{F} . The bracketing entropy is the logarithm of the bracketing number. In the definition of the bracketing number, the upper and lower bounds u and l of the brackets need not belong to \mathcal{F} themselves but are assumed to have finite norms.*

Bracketing number can be regarded as an analog of covering number, describing the geometric complexity of the underlying function space. Although bracketing number of a general function class is difficult to characterize, for some specific function classes, we can easily derive upper bounds for their bracketing number. In particular, we have the following result for Lipschitz functions.

Theorem 5 (Adapted from Theorem 2.7.11 in [37]). *Let $\mathcal{F} = \{f_t : t \in T\}$ be a class of functions. Suppose that for arbitrary $s, t \in T$, we have*

$$|f_s(x) - f_t(x)| \leq d(s, t)F(x), \quad (138)$$

for some metric d on the index set, function F on the sample space, and every x . Then, for any norm $\|\cdot\|$,

$$N_{[\cdot]}(2\varepsilon\|F\|, \mathcal{F}, \|\cdot\|) \leq N(\varepsilon, T, d). \quad (139)$$

Theorem 6 (Adapted from Theorem 2.14.2 in [37]). *For a given norm $\|\cdot\|$, define a bracketing integral of a class of functions \mathcal{F} as*

$$J_{[\cdot]}(\delta, \mathcal{F}, \|\cdot\|) = \int_0^\delta \frac{1}{1 + \log N_{[\cdot]}(\varepsilon\|F\|, \mathcal{F}, \|\cdot\|)} d\varepsilon. \quad (140)$$

Let \mathcal{F} be a class of measurable functions with measurable envelope function F , we have

$$E \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(X_i) - E[f(X)] \leq J_{[\cdot]}(1, \mathcal{F}, L_2(P)) \frac{\|F\|_{L_2(P)}}{\sqrt{n}}, \quad (141)$$

where P is the distribution of X , and the $L_2(P)$ -norm is defined as $\|f\|_{L_2(P)} := \left(\int_{\Omega} f^2(\omega) dP(\omega) \right)^{1/2}$.