

---

# A Convex Relaxation Approach to Generalization Analysis for Parallel Positively Homogeneous Networks

---

Uday Kiran Reddy Tadipatri  
University of Pennsylvania

Benjamin D. Haeffele  
University of Pennsylvania

Joshua Agterberg  
University of Illinois Urbana-Champaign

René Vidal  
University of Pennsylvania

## Abstract

We propose a general framework for deriving generalization bounds for parallel positively homogeneous neural networks—a class of neural networks whose input-output map decomposes as the sum of positively homogeneous maps. Examples of such networks include matrix factorization and sensing, single-layer multi-head attention mechanisms, tensor factorization, deep linear and ReLU networks, and more. Our general framework is based on linking the non-convex empirical risk minimization (ERM) problem to a closely related convex optimization problem over prediction functions, which provides a global, achievable lower-bound to the ERM problem. We exploit this convex lower-bound to perform generalization analysis in the convex space while controlling the discrepancy between the convex model and its non-convex counterpart. We apply our general framework to a wide variety of models ranging from low-rank matrix sensing, to structured matrix sensing, two-layer linear networks, two-layer ReLU networks, and single-layer multi-head attention mechanisms, achieving generalization bounds with a sample complexity that scales almost linearly with the network width.

establishing guaranteed performance of such models—particularly regarding theoretical guarantees on unseen data. This lack of performance guarantees is especially concerning for high-stakes applications such as autonomous vehicles, healthcare, or other high-consequence decision-making systems. To ensure the safe and reliable deployment of deep learning models, it is essential that generalization guarantees be established under reasonable data-generating mechanisms.

**Related work.** There is a broad literature on generalization theory. Classical approaches can be categorized along two separate (but related) lines: (i) data-dependent versus data-independent bounds, and (ii) uniform versus non-uniform concentration. Informally, *data-dependent* bounds take into account explicit data-generating assumptions, whereas data-independent bounds hold *regardless* of the underlying data distribution. Similarly, *uniform* concentration guarantees focus on obtaining concentration inequalities *simultaneously* for all functions in some function class (known as the *hypothesis space*). In contrast, non-uniform concentration inequalities focus on particular functions estimated from the data. Classical approaches marry these two separate types of analyses by introducing measures such as the VC-dimension (Vapnik and Chervonenkis, 1968) or the Rademacher Complexity (Bartlett and Mendelson, 2001). However, these classical measures are often difficult to compute and overly pessimistic, especially when applied to DNNs (Zhang et al., 2021). Consequently, many classical approaches may fail in the modern, more complex DNN setting.

Modern generalization frameworks for DNNs acknowledge that data often comes from structured distributions (e.g., with an intrinsic dimensionality significantly below that of the ambient space) and that optimization algorithms like Stochastic Gradient Descent (SGD) explores only a small portion of the hypothe-

## 1 INTRODUCTION

Despite significant recent advances in the analysis of deep neural networks (DNNs), key gaps persist in es-

sis space (Neyshabur et al., 2017). As a result, the effective hypothesis space is much smaller than what classical bounds account for based on the expressivity of the model alone. Consequently, modern bounds focus on data-dependent, non-uniform approaches. For instance, *margin bounds* (Neyshabur et al., 2018; Golowich et al., 2018; Barron and Klusowski, 2019) provide specific generalization error bounds for DNNs trained to minimize max-margin type loss functions for classification tasks. Another line of research (Dziugaite and Roy, 2017; Arora et al., 2018; Banerjee et al., 2020) exploits the sensitivity of the non-convex landscapes around learned weights; however, this approach requires the estimation of hard quantities like expected sharpness and KL divergence and questions remain regarding the extent to which quantities such as sharpness explain network generalization (Wen et al., 2023; Andriushchenko et al., 2023).

Recent work has observed that optimization methods such as SGD, even without explicit regularization, tend to yield solutions that generalize well, a notion known as *implicit bias* (Gunasekar et al., 2017, 2018a,b; Soudry et al., 2018; Li et al., 2020; HaoChen et al., 2021; Vardi, 2023). This stands in contrast to classical theory, which suggests that explicit regularization is necessary to avoid overfitting. For example, DNNs have been shown to converge toward maximum-margin solutions in classification tasks (Soudry et al., 2018), while solutions in regression tasks often exhibit low-rank structures (Li et al., 2020) that generalize well. Although these analyses provide valuable insights, they are generally limited to specific objectives and types of neural network architectures.

A key challenge in understanding the generalization properties of DNNs is their non-convex landscape. Indeed, convex landscapes are better understood, and numerous generalization bounds have already been derived (Shalev-Shwartz et al., 2009; Lugosi and Neu, 2022). We argue that bridging the gap between non-convex and convex landscapes could provide a pathway to understanding generalization better. Our key contribution is to propose a new generalization analysis framework for DNNs based on linking their non-convex landscape to a convex one. Our framework builds upon Haeffele and Vidal (2017) and Vidal et al. (2022), who connected certain non-convex optimization problems to closely related convex ones. However, their work focuses on characterizing the optimization properties of such problems and does not consider generalization.

**Paper contributions.** In this work, we use the idea of analyzing non-convex problems via a closely related convex problem to derive generalization bounds for a broad family of learning models, which take the form of sums of (*slightly generalized*) positively homogeneous

functions whose parameters are regularized by sums of positively homogeneous functions of the same degree. This allows for a reinterpretation of the (empirical and expected) non-convex optimization problems as closely related to carefully constructed convex problems. We then apply concentration of measure techniques to the convexified version under reasonable data distributions and show that this also implies the concentration of the non-convex problem of interest. More specifically, we extend the finite-dimensional framework of Haeffele and Vidal (2017) and Vidal et al. (2022) to its infinite-dimensional counterpart, which allows us to derive generalization guarantees from a novel viewpoint by exploiting the connection between our problem of interest and a closely related convex problem. We note that other prior work (Bach, 2017) has also considered similar relationships between convex and non-convex problems for establishing generalization results. However, the generalization guarantees in Bach (2017) largely rely on Rademacher complexities, which results in a sample complexity that grows quadratically with the network width. In contrast, we exploit the relationship between the convex and non-convex problems more directly, which allows us to derive bounds with an improved sample complexity.

To be more precise, our main results can be stated informally as follows. Let  $N$  be the number of data points,  $R$  be the number of positively homogeneous functions (or the width of the network) whose predictions are summed together to form the output, and  $\dim(\mathcal{W})$  be the dimension of the parameters in one of the functions. When  $N \gtrsim \tilde{O}(R \times \dim(\mathcal{W}))$ , we show that the generalization error can be bounded with high probability by two terms: the first term, dubbed the optimization error, which vanishes at a globally optimal solution, and the second term, dubbed the statistical error, which depends on the ratio  $\frac{R \times \dim(\mathcal{W})}{N}$ , and hence vanishes only asymptotically.

Our results apply to a wide range of signal processing and DNN problems. The derived bounds achieve near state-of-the-art sample complexity for non-convex low-rank matrix sensing that match the lower bound provided by Candès and Plan (2011) for convex low-rank matrix sensing. By applying these general results to two-layer linear (and ReLU) neural networks with weight decay and multi-head attention models, a key component of transformer architecture (Vaswani et al., 2017), we obtain novel generalization bounds with “tight”<sup>1</sup> sample complexities for both problems.

<sup>1</sup>Our notion of “tight” bounds corresponds to cases where the sample complexity scales linearly or nearly linearly (up to logarithmic factors) with the number of model parameters.

**Outline.** The remainder of this paper is organized as follows. In §2, we formulate the learning problem and introduce our approach. In §3, we explore how learning problems can be bounded via convex surrogates. In §4, we present the statistical bounds through the master theorem that provides generalization error bounds. In §5, we apply the master theorem to various problems in signal processing and DNNs and compare our derived sample complexities with those in the existing literature. The supplementary material contains detailed proofs of the mathematical statements, validations of our framework’s assumptions through simulations, and an additional survey of related works.

**Notation.** For two random variables  $(Z, W)$ , drawn from a joint distribution  $q$ , we define  $\langle Z, W \rangle_q = \mathbb{E}[\langle Z, W \rangle]$ , where the expectation is with respect to the joint probability distribution  $q$ . For a generic function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , we denote  $\|f\|_{Lip}$  as its Lipschitz constant; i.e., the smallest number  $L_f$  such that  $|f(x) - f(y)| \leq L_f \|x - y\|$ . A function  $f$  is said to be integrable with respect to measure  $q$ , i.e.,  $f \in L^2(q)$ , if  $(\int_{x \in \mathcal{X}} \|f(x)\|^2 dq(x))^{1/2} < \infty$ . The inequality  $f(x) \gtrsim g(x)$ , means that there exists a constant  $c > 0$  such that  $f(x) \geq cg(x)$ . We define the ReLU function as  $[x]_+ = \max(x, 0)$ . For the matrix  $U$  the variable  $\mathbf{u}_j$  corresponds to  $j$ th column of  $U$ .

## 2 PROBLEM FORMULATION

Given a realization of a pair  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$  from a distribution  $\mu$  with  $\mathcal{X} \subset \mathbb{R}^{n_x}$ ,  $\mathcal{Y} \subset \mathbb{R}^{n_y}$ , we consider a (non)parametric regression problem of the form  $Y = g(X, \epsilon)$ , where  $\epsilon$  is a source of additional noise (typically independent from  $X$ ). We are interested in approximating  $g$  by the sum of  $r$  prediction functions,  $\phi : \mathcal{W} \times \mathbb{R}^{n_x} \rightarrow \mathbb{R}^{n_y}$ , parameterized by  $W \in \mathcal{W}$ , i.e.,

$$\hat{Y} = \sum_{j=1}^r \phi(W_j)(X) = \Phi_r(\{W_j\})(X). \quad (1)$$

We will additionally refer to  $\phi(W)(X)$  as the *factor map/sub-network* depending on the specific problem.

Our goal is to learn the parameters  $\{W_j\}_{j=1}^r$  that minimize the regularized population risk defined as

$$\begin{aligned} \text{NC}_\mu(\{W_j\}) := & \underbrace{\mathbb{E}_{(X,Y)} [\ell(Y, \Phi_r(\{W_j\}_{j=1}^r)(X))]}_{=:\ell(g, \Phi_r(\{W_j\}_{j=1}^r))_\mu} \\ & + \lambda \Theta_r(\{W_j\}_{j=1}^r), \end{aligned} \quad (2)$$

where  $Y = g(X, \epsilon)$  is the target random variable,  $\ell(\cdot, \cdot)$  is the *loss function*, typically convex in the second argument, and  $\Theta_r(\{W_j\}_{j=1}^r)$  is an *explicit regularization*

function which helps find structured parameters, such as minimum norm or sparse solutions. Specifically, the regularization term  $\Theta_r(\{W_j\}_{j=1}^r)$  is defined as

$$\Theta_r(\{W_j\}_{j=1}^r) := \sum_{j=1}^r \theta(W_j), \quad (3)$$

where  $\theta : \mathcal{W} \rightarrow \mathbb{R}^+$  is a regularization term for each factor map, and  $\lambda \in \mathbb{R}^+$  is a regularization hyperparameter that controls the trade-off between loss reduction and inducing structure.

Notice that we will minimize the population risk  $\text{NC}_\mu(\{W_j\})$  over both  $r$  and  $\{W_j\}_{j=1}^r$ . More explicitly, we will allow for problems where, in addition to optimizing over the model parameters, one also optimizes over the number of prediction functions  $r$  (e.g., the network width) during training. However, our results will also apply to a value of  $r$  that is fixed *a priori*.

Estimating  $\text{NC}_\mu(\{W_j\})$  directly is challenging due to (i) the lack of access to the distribution  $\mu$ , (ii) the fact that  $(\{W_j\})$  (and potentially the number  $r$ ) are random variables dependent on the training data  $\{(X_i, Y_i)\}$ , and (iii) the non-linearity and potential non-convexity of  $\text{NC}_\mu$ . We address the first point (as is standard) via empirical minimization of  $\text{NC}_\mu(\cdot)$  using the *empirical risk* (or *training error*) defined via:

$$\begin{aligned} \text{NC}_{\mu_N}(\{W_j\}) := & \frac{1}{N} \sum_{i=1}^N \ell(Y_i, \Phi_r(\{W_j\}_{j=1}^r)(X_i)) \\ & \underbrace{=:\ell(g, \Phi_r(\{W_j\}_{j=1}^r))_{\mu_N}}_{+ \lambda \Theta_r(\{W_j\}_{j=1}^r)}, \end{aligned} \quad (4)$$

where  $\mu_N$  denotes the empirical distribution of the samples  $\{X_i, Y_i\}_{i=1}^N$ . We define empirical risk minimization (ERM) via the arg min of  $\text{NC}_{\mu_N}(\{W_j\})$ . For concreteness, recall we also allow for the minimization over  $r$  (provided  $r$  is bounded above by some quantity independent of the data), though our results hold for any fixed  $r$ .

Note that if we minimize the objective  $\text{NC}_{\mu_N}(\cdot)$ , there is no guarantee that we will also minimize  $\text{NC}_\mu(\cdot)$ . This discrepancy is quantified by the *Generalization Error*:

$$\begin{aligned} & |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \\ & = |\ell(g, \Phi_r(\{W_j\}_{j=1}^r))_\mu - \ell(g, \Phi_r(\{W_j\}_{j=1}^r))_{\mu_N}|. \end{aligned} \quad (5)$$

Note that the regularization terms containing  $\Theta_r$  are the same between the two objectives, giving the typical difference between the empirical and population losses.

In this work, we compute an upper bound for the generalization error at any stationary point of the empirical problem,  $\text{NC}_{\mu_N}(\{W_j\})$ , under certain technical assumptions. To build our main results, we relate these

<sup>2</sup>We occasionally notate  $\{W_j\}_{i=1}^r$  as  $\{W_j\}$  for brevity of notation, but the dependence on  $r$  is always implied.

non-convex objectives  $\text{NC}_\mu(\{W_j\})$  and  $\text{NC}_{\mu_N}(\{W_j\})$  to closely related *convex* objectives in the prediction space, respectively  $\text{C}_\mu(f_\mu)$  and  $\text{C}_{\mu_N}(f_{\mu_N})$ , whose definitions will be introduced in §3. This allows us to decompose the generalization error in (5) as:

$$\begin{aligned} \text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\}) &= \underbrace{\left[ \text{NC}_\mu(\{W_j\}) - \text{C}_\mu(f_\mu) \right]}_{\text{Population Gap}} \\ &\quad - \underbrace{\left[ \text{NC}_{\mu_N}(\{W_j\}) - \text{C}_{\mu_N}(f_{\mu_N}) \right]}_{\text{Empirical Gap}} + \underbrace{\left[ \text{C}_\mu(f_\mu) - \text{C}_{\mu_N}(f_{\mu_N}) \right]}_{\text{Convex Generalization Gap}}. \end{aligned} \quad (6)$$

Our Theorem 1 bounds the *Empirical Gap* and the *Population Gap*. With these bounds, we then apply concentration techniques to bound the *Convex Generalization Gap* and obtain our main Theorem 2 which gives bounds for the generalization error in (5).

### 3 CONVEX BOUNDS FOR LEARNING

In this section, we present bounds for the *Empirical Gap* and *Population Gap* through Theorem 1, linking our learning problem of interest to functions that are convex in the space of prediction functions. To begin, we state several requirements for our framework.

**Assumption 1** (Regularization). *The regularization function  $\theta$  is positive semidefinite; i.e.,  $\theta(0) = 0$  and  $\theta(W) \geq 0, \forall W \in \mathcal{W}$ .*

This is a mild assumption; it only ensures we do not impose negative regularization on the parameters  $\{W_j\}$ . Our next assumption is our main functional assumption on  $\phi$  and  $\theta$ .

**Assumption 2** (Balanced Homogeneity of  $\phi$  and  $\theta$ ). *The factor map  $\phi$  and the regularization map  $\theta$  can be scaled equally by non-negative scaling of (a subset of) the parameters. Formally, we assume that there exists sub-parameter spaces  $(\mathcal{K}, \mathcal{H})$  from the parameter space  $\mathcal{W}$  such that  $\mathcal{K} \times \mathcal{H} = \mathcal{W}$ ,  $\forall (\mathbf{k}, \mathbf{h}) \in (\mathcal{K}, \mathcal{H})$ , and  $\beta \geq 0$  we have  $\phi((\beta\mathbf{k}, \mathbf{h})) = \beta^p \phi((\mathbf{k}, \mathbf{h}))$  and  $\theta((\beta\mathbf{k}, \mathbf{h})) = \beta^p \theta((\mathbf{k}, \mathbf{h}))$  for some  $p > 0$ . Further, we assume that for bounded input  $X$  the set  $\{\phi(W)(X) : \forall W \in \mathcal{W} \text{ s.t. } \theta(W) \leq 1\}$  is bounded.*

This is a slight generalization of positive homogeneity, which only requires positive homogeneity in a subset of parameters, provided the image of the factor map for parameters with  $\theta(W) \leq 1$  is bounded<sup>3</sup>.

<sup>3</sup>For example,  $\phi(v)(X)$  can take the form  $v^a g(X)$  for  $a \geq 0$ , and  $\theta(v) = |v|^a$ , where  $g : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$  is some fixed function. More generally, we can choose  $\phi(v_1, v_2) = v_1^a g_{v_2}(X)$  and  $\theta(v_1, v_2) = |v_1|^a + \delta_{\mathcal{V}_2}(v_2)$ , where  $g_{v_2} : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$  is a function parameterized by  $v_2 \in \mathcal{V}_2$  and has a bounded range for bounded inputs.

Our next assumption concerns the loss function  $\ell$ .

**Assumption 3** (Convex Loss). *The loss  $\ell(Y, \hat{Y})$  is second-order differentiable (written  $\ell \in \mathcal{C}^2$ ), and  $L$ -smooth w.r.t.  $\hat{Y}$ , i.e., for any  $Y, \hat{Y} \in \mathbb{R}^{n_Y}$*

$$0 \preceq \nabla_{\hat{Y}}^2 \ell(Y, \hat{Y}) \preceq L I_{n_Y}. \quad (7)$$

*Additionally, the gradient of the loss is bi-Lipschitz smooth; that is, for all  $Y_1, Y_2, \hat{Y}_1, \hat{Y}_2 \in \mathbb{R}^{n_Y}$*

$$\begin{aligned} \|\nabla_{\hat{Y}} \ell(Y_2, \hat{Y}_2) - \nabla_{\hat{Y}} \ell(Y_1, \hat{Y}_1)\| &\leq L \left[ \|Y_2 - Y_1\|_2 \right. \\ &\quad \left. + \|\hat{Y}_2 - \hat{Y}_1\|_2 \right], \end{aligned} \quad (8)$$

*and the loss is constant if both the arguments are the same, i.e., for all  $Y_1, Y_2 \in \mathbb{R}^{n_Y}$ ,  $\ell(Y_1, Y_1) = \ell(Y_2, Y_2)$ .*

This ensures that the loss function is convex and smooth. Furthermore, if the loss is  $\alpha$ -strongly convex, i.e.,  $0 \prec \alpha I_{n_Y} \preceq \nabla_{\hat{Y}}^2 \ell(Y, \hat{Y}) \preceq L I_{n_Y}$  we have derived tighter results (see the Appendix).

We define the *induced regularization* function as

$$\begin{aligned} \Omega(f) &:= \inf_{r, \{W_j\}} \Theta_r(\{W_j\}) \\ \text{s.t. } f(X) &= \Phi_r(\{W_j\})(X); \quad \forall X \in \mathcal{X}, \end{aligned} \quad (9)$$

with the function taking value infinity if  $f(X)$  cannot be realized for some choice of the parameters  $(r, \{W_j\}_{j=1}^r)$ . Using similar arguments as in Haeffele and Vidal (2015) it can be shown that under assumptions 1–2, the function  $\Omega(f)$  is convex in the space of prediction functions; see Proposition 1 in the Appendix. Moreover, by Assumption 3, the loss function is convex with respect to the model predictions, which allows us to define the following two *convex* optimization problems over the space of prediction functions:

$$\text{C}_\mu(f) := \mathbb{E}_{(X, Y)}[\ell(Y, f(X))] + \lambda \Omega(f), \quad (10)$$

where  $f \in L^2(\mu)$ , and

$$\text{C}_{\mu_N}(f) := \frac{1}{N} \sum_{i=1}^N \ell(Y_i, f(X_i)) + \lambda \Omega(f), \quad (11)$$

where  $f \in L^2(\mu_N)$ .

From the definition of  $\Omega(f)$  we have that  $\text{C}_\mu$  and  $\text{C}_{\mu_N}$  are always lower bounds of  $\text{NC}_\mu$  and  $\text{NC}_{\mu_N}$ , respectively, for any  $(f, \{W_j\})$  such that  $f(X) = \Phi_r(\{W_j\})(X)$ , which becomes a tight bound for any parametrization  $(\{W_j\})$  of  $f$  which achieves the infimum. As a result, we can relate solutions of the non-convex problems to the corresponding convex problem via tools from convex analysis, as we establish in the following result.

**Theorem 1** (Convex Bounds for Learning). *Under assumptions 1–3, let  $f_{\mu_N}^*$  (or  $f_\mu^*$ ) be the global minimizer for  $C_{\mu_N}(\cdot)$  (or  $C_\mu(\cdot)$ ). For any stationary points  $(r, \{W_j\})$  of the function  $NC_{\mu_N}(\cdot)$  and any  $f \in L^2(\mu) \cap L^2(\mu_N)$  the following are true:*

1. *Empirical optimality gap:*

$$C_{\mu_N}(f_{\mu_N}^*) \leq NC_{\mu_N}(\{W_j\}) \leq C_{\mu_N}(f) + \lambda \Omega(f) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right], \quad (12)$$

2. *Population optimality gap:*

$$C_\mu(f_\mu^*) \leq NC_\mu(\{W_j\}) \leq C_\mu(f) + \lambda \Omega(f) \left[ \Omega_\mu^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] + \left[ \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N} \right], \quad (13)$$

where  $\Omega_q^\circ(\cdot)$  is referred to as polar in the measure  $q$  defined as

$$\Omega_q^\circ(g) := \sup_{\theta(W) \leq 1} \langle g, \phi(W) \rangle_q. \quad (14)$$

Readers are referred to Appendix A for the proof with extensions to strongly convex functions.

The population optimality gap is obtained by an infinite-dimensional extension of Proposition 3 in Haeffele and Vidal (2020). The additional term in the population optimality gap (13) arises from the fact that the stationary points of ERM,  $NC_{\mu_N}(\cdot)$ , are not necessarily the same as those of  $NC_\mu(\cdot)$ .

From equation (5) the goal is to bound the difference between the original non-convex formulations  $NC_{\mu_N}$  and  $NC_\mu$ . By Theorem 1, we established the optimality gaps for both empirical and population non-convex optimization problems, and by computing the difference between equation (12) and (13), with algebraic manipulation we arrive at the following quantities:

- **Convex Generalization Gap:** The convex generalization gap is defined as  $|C_\mu(f) - C_{\mu_N}(f)|$ .
- **Polar Gap:** By virtue of the fact that the loss functions each contain the respective polars, we define the Polar Gap as the quantity  $|\Omega_{\mu_N}^\circ(\nabla_{\hat{Y}} \ell(g, f)) - \Omega_\mu^\circ(\nabla_{\hat{Y}} \ell(g, f))|$ .
- **Equilibria Gap:** We define the Equilibria Gap via  $\left| \langle \nabla_{\hat{Y}} \ell(g, f), f \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f), f \rangle_\mu \right|$ .

- **Norm Gap:** The final remaining quantity is defined via  $\left| \|f_\mu^* - f\|_{\mu_N}^2 - \|f_\mu^* - f\|_\mu^2 \right|$ . This quantity applies only to strongly convex functions (see the Appendix).

A major technical contribution of this paper is to demonstrate that each of these quantities uniformly concentrates at a rate equal to or smaller than the “statistical error” under certain realistic assumptions that are discussed in §4. The only remaining term from Theorem 1 is the quantity  $\Omega(f_\mu)[\Omega_{\mu_N}^\circ(\cdot) - 1]$ , which bounds the sub-optimality (in objective value) of the current stationary point for the empirical optimization problem. This term approaches zero at the global optimum of  $NC_{\mu_N}$  (see §A in the Appendix).

## 4 STATISTICAL BOUNDS

In Theorem 1, we established bounds for the *Empirical Gap* and *Population Gap*. Building on these results, we identified key quantities such as the *Convex Generalization Gap*, *Polar Gap*, *Equilibrium Gap*, and *Norm Gap*, all of which can be controlled under certain general conditions (Assumptions 1–6, along with Assumption 7’ from the Appendix) that we state momentarily. In this section, we present Theorem 2, which consolidates these bounds to derive our main generalization error bound. For clarity and to minimize technical complexity, we present Theorem 2 with Assumption 7, which a stronger version of Assumption 7’.

To begin, we state our additional assumptions. We assume that  $\phi$  is Lipschitz.

**Assumption 4** (Lipschitz Continuity of  $\phi$ ). *Let  $\mathcal{B}$  be some compact subset of  $\mathcal{W}$ , and denote*

$$\mathcal{F}_\theta := \{W : \theta(W) \leq 1\} \cap \mathcal{B} \subseteq \mathbb{B}(r_\theta), \quad (15)$$

where  $\mathbb{B}(r_\theta)$  is the  $L_2$  ball with radius  $r_\theta$ .<sup>4</sup> The factor map  $\phi$  is Lipschitz continuous with respect to inputs for any choice of parameters  $W \in \mathcal{F}_\theta$ , i.e.,

$$L_\phi := \sup_{W \in \mathcal{F}_\theta} \|\phi(W)\|_{Lip} < \infty. \quad (16)$$

Our next assumption imposes tail conditions on the random variables  $(X, Y)$ .

**Assumption 5** (Data Model). *The input data  $X \in \mathbb{R}^{n \times x}$  is drawn from the 1-Lipschitz concentrated sub-Gaussian distribution with a proxy variance  $\sigma_X^2/n_X$ ;*

<sup>4</sup>The radius  $r_\theta$  can depend on the dimension of  $W$ . For instance, suppose  $W \in \mathbb{R}^n$  and  $\theta(W) = \|W\|_1$ , as  $\|W\|_1 \leq \sqrt{n}\|W\|_2$ , then  $r_\theta$  must be at least  $\sqrt{n}$ . On another instance, suppose  $W = (\mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^n)$ , and  $\theta(W) = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2$ ; this requires  $r_\theta$  to be at least  $1/2$ .

i.e., for any 1-Lipschitz continuous function,  $h : \mathbb{R}^{n_X} \rightarrow \mathbb{R}$  there exists  $c > 0$  such that

$$P(|h(X) - \mathbb{E}_X[h(X)]| \geq \epsilon) \leq c \exp\left(-\frac{n_X \epsilon^2}{2\sigma_X^2}\right). \quad (17)$$

The target function  $Y$  takes the form  $Y = g(X, \epsilon)$ , where  $g \in L^2(\mu)$  is bi-Lipschitz in  $X$  and  $\epsilon$ ; that is,

$$\|g(X_2, \epsilon_2) - g(X_1, \epsilon_1)\|_2 \leq \|g\|_{Lip} \left[ \|X_2 - X_1\|_2 + \|\epsilon_2 - \epsilon_1\|_2 \right], \quad (18)$$

and  $\epsilon \sim \mathcal{N}(0, (\sigma_{Y|X}^2/n_E)I)$  in  $\mathbb{R}^{n_E}$ .

We note that the above assumption is mild. While extending our framework to heavy-tailed distributions are likely possible; it would require a more intricate analysis and may result in worse error rates and larger sample complexities.

Our next assumption concerns the possible functions learned via empirical risk minimization.

**Assumption 6** (Hypothesis class). *Stationary points of  $\text{NC}_{\mu_N}(\cdot)$  have bounded regularization and bounded width,  $r \leq R$ , almost surely. The input-output map,  $\Phi_r(\{W_j\})$  has Lipschitz constant at most  $\gamma$ , and the parameters are bounded. Let  $\mathcal{B}_R \subseteq \mathcal{W}^R$  be some compact set; then the hypothesis class is defined as*

$$\mathcal{F}_{\mathcal{W}} := \left\{ \{W_j\}_{j=1}^r : \|\Phi_r(\{W_j\})\|_{Lip} \leq \gamma \right\} \cap \mathcal{B}_R. \quad (19)$$

In words, the set of maps learned through ERM are essentially Lipschitz in the parameters  $\{W_j\}$ , and, furthermore, the  $\{W_j\}$  are bounded (almost surely). Moreover, the assumption that  $r \leq R$  ensures that at most  $R$  individual functions  $\{W_j\}$  are needed, which implicitly imposes a “low-complexity” constraint on the learned function. Finally, note that we assume that  $\gamma$  does not depend on the width of the network. In practice, our empirical observations show that the Lipschitz constant does not increase with width, making it a realistic assumption. For further details, refer to the numerical simulations in §E of Appendix.

Our general master theorem, Theorem 4 in the Appendix, requires only Assumptions 1–6 and 7’ (in the Appendix). For the sake of notational brevity, we state our main results with the slightly stronger Assumption 7 instead of Assumption 7’.

**Assumption 7** (Boundedness). *For all  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ , and  $\{W_j\} \in \mathcal{F}_{\mathcal{W}}$ , the predictions, and gradients are bounded; i.e.,*

$$\|\Phi_r(\{W_j\})(X)\| \leq B_{\Phi}, \quad \|\nabla_{\hat{Y}} \ell(Y, \Phi_r(\{W_j\})(X))\| \leq B_{\ell}. \quad (20)$$

Further, for any  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ , for any  $\{W_j\}, \{\tilde{W}_j\} \in \mathcal{F}_{\mathcal{W}}$ ,  $W, \tilde{W} \in \mathcal{F}_{\theta}$ , the network,  $\phi$  and  $\Phi_r$ , are Lipschitz in the parameters; i.e.,

$$\|\Phi_r(\{W_j\})(X) - \Phi_r(\{\tilde{W}_j\})(X)\|_2 \leq \tilde{L}_{\Phi} \max_j \|W_j - \tilde{W}_j\|_2, \text{ and} \quad (21)$$

$$\|\phi(W)(X) - \phi(\tilde{W})(X)\|_2 \leq \tilde{L}_{\phi} \|W - \tilde{W}\|_2. \quad (22)$$

Assumption 7 ensures that predictions and its gradients are bounded while the network being Lipschitz continuous on the parameter space for any inputs. Assumption 7 implicitly indicates that either the data points are uniformly bounded or the search space for the parameters is of small dimension, which can restrict the potential applications. However, as we demonstrate in the more general version (Theorem 4) in the Appendix, it suffices that the conditions above hold only for some convex set  $\mathcal{C}$ , though this extension requires significantly more notation and discussion, so we do not include it here.

**Theorem 2** (Master Theorem). *Suppose Assumptions 1–7 hold. Let  $\delta \in (0, 1]$  be fixed, and let  $f_{\mu}^*$  be the global optimum of  $\mathcal{C}_{\mu}$ . Suppose that  $\gamma \geq \Omega(f_{\mu}^*) \tilde{L}_{\phi}$ , and define*

$$\epsilon_1 = 16\gamma^2 \sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{Lip}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}; \quad (23)$$

$$\epsilon_2 = 4\tilde{L}_{\Phi} B_{\Phi} \max \left\{ 1, 2L + 2B_{\ell}/B_{\Phi}, 8\Omega(f_{\mu}^*)(B_{\ell} \tilde{L}_{\phi})/(\tilde{L}_{\Phi} B_{\Phi}), 8L\Omega(f_{\mu}^*) \right\}. \quad (24)$$

Let  $\{W_j\}$  denote any stationary point of  $\text{NC}_{\mu_N}(\cdot)$ . Then with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned} \frac{1}{n_Y} |\text{NC}_{\mu}(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| &\lesssim \quad (25) \\ &\underbrace{\frac{\lambda}{n_Y} \Omega(f_{\mu}^*) \left[ \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\}))_{\mu_N} \right) - 1 \right]}_{\text{Optimization Error}} \\ &+ \epsilon_1 \underbrace{\sqrt{\frac{R \cdot \dim(\mathcal{W}) \log\left(\frac{\gamma \epsilon_2 r \theta}{L_{\phi}}\right) \log(N) + \log\left(\frac{1}{\delta}\right)}{N}}}_{\text{Statistical Error}} \end{aligned}$$

**Remarks:** The generalization error is upper bounded by two terms:

- the *Optimization Error*, which quantifies the distance to the globally optimal solution, and
- the *Statistical Error*, or the intrinsic error that depends on the sample complexity and the noise.

The optimization error diminishes as we approach a global optimum of the ERM problem  $\text{NC}_{\mu_N}$  and vanishes at a global optimum, whereas the statistical error diminishes as the sample size increases relative to

the intrinsic dimension, i.e., when  $N \gtrsim R \times \dim(\mathcal{W})$  (ignoring logarithmic factors). By a naive counting argument, there are  $R \times \dim(\mathcal{W})$  many parameters in the underlying network. As we will see in subsequent sections, this sample complexity turns out to be optimal or nearly optimal for a number of reasonable statistical settings. The implicit constants appearing in the result are universal and are not problem dependent.

## 5 APPLICATIONS

In this section, we present applications of the Theorem 2 for low-rank matrix sensing, two-layer ReLU neural networks, and single-layer multi-head attention. To apply Theorem 2, we must compute the problem-specific quantities  $\Omega(f_\mu^*)$ ,  $\Omega_{\mu_N}^\circ(\cdot)$ ,  $L$ ,  $\|g\|_{\text{Lip}}$ ,  $\sigma_X$ ,  $\sigma_{Y|X}$ ,  $\epsilon_1$ ,  $\epsilon_2$ ,  $r_\theta$ ,  $\gamma$ ,  $L_\phi$ . For each application, we have estimated these quantities, with further details provided in the proofs located in Appendix C.1, C.4, and C.5, respectively. We summarize and compare the obtained sample complexities for the various applications with their state-of-the-art bounds in Table 1. The additional applications to structured matrix sensing and two-layer linear neural networks can be found in Appendix C.2 and Appendix C.3, respectively.

**Low-rank matrix sensing:** We first consider low-rank matrix sensing (Candès and Plan, 2011), which is a well-studied problem in the signal processing and statistics literature. Given a few linear measurements of an unknown low-rank matrix, the goal is to estimate the low-rank matrix in the presence of noise. One potential strategy is to define a convex program via nuclear-norm regularization (Candès and Recht, 2009). While recovery guarantees for this convex program are well-studied, solving it is a computationally intensive procedure involving computing a full singular value decomposition at each iteration. To address this issue, several authors have considered a non-convex variant that reparameterizes the low-rank matrix into its underlying left and right factors, which is known as the Burer-Monteiro factorization (Burer and Monteiro, 2003). While the new optimization problem runs faster in practice, it is also non-convex, and its properties can be difficult to analyze theoretically. Corollary (1) provides the bounds on the generalization error for this non-convex program.

**Corollary 1** (Low-Rank Matrix Sensing). *Consider the true model for  $(X, y)$ , where  $X \in \mathbb{R}^{m \times n}$  is a random matrix with i.i.d. entries  $X_{lk} \sim \mathcal{N}(0, \frac{1}{mn})$  and  $y = \langle M^*, X \rangle + \epsilon$ , where  $M^* \in \mathbb{R}^{m \times n}$  and  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  is independent from  $X$ . For all  $i \in [N]$ , let  $(X_i, y_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{y} = \langle UV^T, X \rangle$ , where  $U \in \mathbb{R}^{m \times R}$  and  $V \in \mathbb{R}^{n \times R}$ . Let  $\delta \in (0, 1]$  be fixed. Define the non-*

*convex problem*

$$\text{NC}_{\mu_N}^{\text{MS}}((U, V)) := \frac{1}{2N} \sum_{i=1}^N (y_i - \langle UV^T, X_i \rangle)^2 + \lambda \sum_{j=1}^R \|\mathbf{u}_j\|_2 \|\mathbf{v}_j\|_2, \quad (26)$$

and define  $\text{NC}_\mu^{\text{MS}}((U, V))$  similarly with the sum over  $i$  replaced by expectation taken over  $(X, y)$ .

Let  $(\hat{U}, \hat{V})$  be a stationary point of  $\text{NC}_{\mu_N}^{\text{MS}}(\cdot)$ . Suppose there exists  $C_{UV}, B_u, B_v > 0$  such that  $\|\hat{U}\hat{V}^T\|_2 \leq C_{UV}\|M^*\|_*$ , and for all  $j \in [R]$ ,  $\|\hat{\mathbf{u}}_j\|_2 \leq B_u$ ,  $\|\hat{\mathbf{v}}_j\|_2 \leq B_v$ . Then with probability at least  $1 - \delta$ , it holds that

$$\left| \text{NC}_\mu^{\text{MS}}((\hat{U}, \hat{V})) - \text{NC}_{\mu_N}^{\text{MS}}((\hat{U}, \hat{V})) \right| \lesssim \left\| M^* \right\|_* \left[ \left\| \frac{1}{N} \sum_{i=1}^N (y_i - \langle \hat{U}\hat{V}^T, X_i \rangle) X_i \right\|_2 - \lambda \right] + C_{UV}^2 \|M^*\|_*^2 \times \sqrt{\frac{R \log(R(C_{UV} + B_u B_v)) (m+n) \log(N) + \log(1/\delta)}{N}}. \quad (27)$$

**Remarks:** Observe that at a global minimum, the right side tends to zero when  $R(m+n)/N \rightarrow 0$ , ignoring logarithmic terms. Existing literature on non-convex noisy low-rank matrix sensing typically requires knowledge of true rank( $M^*$ ) =  $R^*$ , and the state-of-the-art sample complexity for this setting is of order  $R^*(m+n)$  in the un-regularized setting (Stöger and Zhu, 2024). In contrast, Corollary 1 does not require knowledge of the true rank. However, if the estimated rank  $R$  is too small ( $R < R^*$ ), then the optimization error still persists. In contrast, if ( $R \geq R^*$ ) then optimization error can vanish subject to the ability of the algorithm utilized to reach stationary points, Haeffele and Vidal (2015) provides such guarantees.

**Two-layer ReLU Networks:** Next, we move on to two-layer ReLU networks, which introduce an additional nonlinearity with respect to the inputs. ReLU networks are widely used and proven to be universal approximators (Huang, 2020). Prior work on generalization analysis for ReLU networks is based on classical measures, such as Rademacher complexity (Bartlett et al., 2019). The following result circumvents the difficulty in the estimate of such classical measures.

**Corollary 2** (Two-Layer ReLU Neural Network). *Consider the true model for  $(\mathbf{x}, \mathbf{y})$ , where  $\mathbf{x} \sim \mathcal{N}(0, (1/n)I_n) \in \mathbb{R}^n$ ,  $\mathbf{y} = U^*[V^{*T}\mathbf{x}]_+ + \epsilon$ , where  $U^* \in \mathbb{R}^{m \times R^*}$ ,  $V^* \in \mathbb{R}^{n \times R^*}$ , and  $\epsilon \sim \mathcal{N}(0, (\sigma^2/m)I_m) \in \mathbb{R}^m$  is independent from  $\mathbf{x}$ . For all  $i \in [N]$ , let  $(\mathbf{x}_i, \mathbf{y}_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{\mathbf{y}} = U[V^T\mathbf{x}]_+$ , where  $U \in \mathbb{R}^{m \times R}$ ,  $V \in \mathbb{R}^{n \times R}$ .*

Table 1: Comparisons with the state-of-the-art sample complexities.  $N$  represents the number of data points.

Application	Our work, $N \gtrsim$	State-of-the-art, $N \gtrsim$
Low rank matrix sensing	$\tilde{\mathcal{O}}(R(m+n))$	$R^*(m+n)$ , (Stöger and Zhu, 2024) (no regularization)
Structured matrix sensing		—
2-Layer linear NN		$R(m+n)$ (Kakade et al., 2008) (bounded data-points)
2-Layer ReLU NN		$R(m+n) \log(R(m+n))$ , (Bartlett et al., 2019)
Multi-head attention		$R(m+n)$ , (Trauger and Tewari, 2024) (bounded data-points)

Let  $\delta \in (0, 1]$  be fixed. Define the non-convex problem

$$\begin{aligned} \text{NC}_{\mu_N}^{\text{ReLU}}((U, V)) &:= \frac{1}{2N} \sum_{i=1}^N \|\mathbf{y}_i - U[V^T \mathbf{x}_i]_+\|_2^2 \\ &\quad + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2), \end{aligned} \quad (28)$$

and define  $\text{NC}_{\mu_N}^{\text{ReLU}}((U, V))$  similarly with the sum over  $i$  replaced by expectation taken over  $(\mathbf{x}, \mathbf{y})$ .

Let  $(\hat{U}, \hat{V})$  be a stationary point of  $\text{NC}_{\mu_N}^{\text{ReLU}}(\cdot)$ . Suppose there exists  $C_{UV}, B_u, B_v > 0$  such that  $\|\hat{U}\hat{V}^T\|_2 \leq C_{UV} [\|U^*\|_F^2 + \|V^*\|_F^2]$ , and for all  $j \in [R]$ ,  $\|\hat{\mathbf{u}}_j\|_2 \leq B_u$ ,  $\|\hat{\mathbf{v}}_j\|_2 \leq B_v$ . Then with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned} \frac{1}{m} \left| \text{NC}_{\mu_N}^{\text{ReLU}}((\hat{U}, \hat{V})) - \text{NC}_{\mu_N}^{\text{ReLU}}((U, V)) \right| &\lesssim \\ \frac{1}{2m} [\|U^*\|_F^2 + \|V^*\|_F^2] &\left[ \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2 \|\mathbf{x}_i\|_2 - \lambda \right] \\ + C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2] \times & \\ \left[ \frac{R(m+n) \log(R(m+n)(C_{UV} + B_u^2 + B_v^2)) \log(N)}{N} \right. & \\ \left. + \frac{\log(1/\delta)}{N} \right]^{1/2}. & \end{aligned} \quad (29)$$

**Remarks:** Analogous to matrix sensing, when  $R(m+n)/N \rightarrow 0$ , the right side tends to zero at global optimality (ignoring logarithmic terms). Furthermore, Corollary 2 recovers the state-of-the-art result by Bartlett et al. (2019).

**Transformers:** Finally, we move on to our last application (though of course, the applications are in fact myriad in principle) to a single layer multi-head attention, which are backbones for transformer-style architecture (Vaswani et al., 2017). In practice, transformers are shown to have remarkable generalization capabilities (Zhou et al., 2024). However, there is a lack of intensive theoretical analysis for this architecture. Few attempts on estimating the capacities of the attention mechanisms have been made in Edelman et al. (2022) and Trauger and Tewari (2024), among others. For our analysis, we consider the case where the output of the

model is one particular token within the input (e.g., transformers use a dedicated class token for the output initialized as a constant vector). The output for one attention head is modeled as  $VX\sigma((KX)^\top Q\mathbf{x}_{out})$  where  $\mathbf{x}_{out}$  is the column of  $X$  corresponding to the transformer output. We then reparameterize  $K^\top Q\mathbf{x}_{out} = \mathbf{z}$  and present the following result.

**Corollary 3 (Transformers).** Consider the true model for  $(X, \mathbf{y})$ , where  $X \in \mathbb{R}^{n \times T}$  is a random matrix with i.i.d. entries  $X_{lk} \sim \mathcal{N}(0, 1/(nT))$  and  $\mathbf{y} = A^*X\mathbf{b}^* + \epsilon$ , where  $A^* \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b}^* \in \mathbb{S}^{T-1}$  and  $\epsilon \sim \mathcal{N}(0, (\sigma^2/m)I_m)$  is independent from  $X$ . For all  $i \in [N]$ , let  $(X_i, \mathbf{y}_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{\mathbf{y}} = \sum_{j=1}^R V_j X \sigma(X^T \mathbf{z}_j)$ ,  $V_j \in \mathbb{R}^n$ ,  $\mathbf{z}_j \in \mathbb{R}^n$ . Let  $\delta \in (0, 1]$  be fixed. Define the non-convex problem

$$\begin{aligned} \text{NC}_{\mu_N}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\}) &:= \frac{1}{2N} \sum_{i=1}^N \|\mathbf{y}_i - \sum_{j=1}^R V_j X_i \sigma_t(X_i^T \mathbf{z}_j)\|_2^2 \\ &\quad + \lambda \sum_{j=1}^R [\|V_j\|_F + \delta_{\{\mathbf{z}: \|\mathbf{z}\|_2 \leq 1\}}(\mathbf{z}_j)], \end{aligned} \quad (30)$$

where,  $\sigma_t(\cdot)$  is softmax function with temperature  $t$ , for  $k \in [T]$  defined  $\sigma_t(\mathbf{u})_k := \exp(tu_k) / \sum_{l=1}^T \exp(tu_l)$  and define  $\text{NC}_{\mu_N}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\})$  similarly with the sum over  $i$  replaced by expectation taken over  $(X, \mathbf{y})$ .

Let  $\{(\hat{V}_j, \hat{\mathbf{z}}_j)\}$  be a stationary point of  $\text{NC}_{\mu_N}^{\text{TF}}(\cdot)$ . Suppose there exists  $C_V, B_V > 0$  such that  $\sum_{j=1}^R \|\hat{V}_j\|_F \leq C_V \|A^*\|_F$ , and for all  $j \in [R]$ ,  $\|\hat{V}_j\|_F \leq B_V$ . Then with probability at least  $1 - \delta$  it holds that

$$\begin{aligned} \frac{1}{m} \left| \text{NC}_{\mu_N}^{\text{TF}}(\{(\hat{V}_j, \hat{\mathbf{z}}_j)\}) - \text{NC}_{\mu_N}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\}) \right| &\lesssim \\ \frac{1}{2m} \|A^*\|_F &\left[ \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2 \|\mathbf{x}_i\|_2 - \lambda \right] \\ + C_V^2 \|A^*\|_F^2 \times & \\ \sqrt{\frac{R(m+n) \log(R(m+n)(C_V + B_V)) \log(N) + \log(1/\delta)}{N}}. & \end{aligned} \quad (31)$$

**Remarks:** The dependence on  $\mathbf{b}^*$  is not explicitly reflected in Equation (31) because the ground truth



model is bilinear. Consequently, assuming  $\mathbf{b}^*$  is unit-norm without loss of generality, as its norm can be absorbed into  $A^*$ . Thus, the dependence on  $\mathbf{b}^*$  is implicitly captured by the norm of  $A^*$  in Equation (31). As in the previous two applications, we can achieve consistency at global optimality when  $N \gtrsim R(m+n)$ , ignoring logarithmic terms. Note that the sample complexity has no dependency on the number of tokens,  $T$ , which suggests an explanation for the success behind the prediction capabilities of transformers for longer length inputs (Zhou et al., 2024). Our sample complexity matches the state-of-the-art bounds on the transformers by Trauger and Tewari (2024).

## 6 CONCLUSIONS

In this work, we provide generalization bounds for non-convex problems of the form of sums of (*slightly generalized*) positively homogeneous functions with a general objective. Our bounds provide sample complexities that are near-optimal and applicable to various problems, such as low-rank matrix sensing, two-layer neural networks, and single-layer multi-head attention. The sample complexity of our bounds grows almost linear with the total number of parameters in the model, and for matrix sensing, this sample complexity is optimal, as demonstrated in Candès and Plan (2011). Our proofs are based on analyzing closely related convex programs in the prediction space; this perspective enabled us to provide near-optimal sample complexities due to existing results on generalization properties for convex functions. In future work, it would be interesting to sharpen the dependence of our bounds on all the relevant parameters and apply our techniques to other machine learning problems.

## Acknowledgments

UKRT gratefully acknowledges Pratik Chaudhari, Hancheng Min, Kyle Poe and Ziqing Xu for their valuable discussions and constructive feedback. His research was supported by the Leggett Family Fellowship and the Dean’s Fellowship programs. Other authors acknowledge the support of the Research Collaboration on the Mathematical and Scientific Foundations of Deep Learning (NSF grant 2031985 and Simons Foundation grant 814201).

## References

Allen-Zhu, Z., Li, Y., and Liang, Y. (2019). Learning and generalization in overparameterized neural networks, going beyond two layers. *Advances in neural information processing systems*, 32.

Andriushchenko, M., Croce, F., Müller, M., Hein, M., and Flammarion, N. (2023). A modern look at the

relationship between sharpness and generalization. *International Conference on Learning Representations (ICLR)*.

- Arora, S., Cohen, N., Golowich, N., and Hu, W. (2019). A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*.
- Arora, S., Ge, R., Neyshabur, B., and Zhang, Y. (2018). Stronger generalization bounds for deep nets via a compression approach. In *International conference on machine learning*, pages 254–263. PMLR.
- Bach, F. (2013). Convex relaxations of structured matrix factorizations. *arXiv preprint arXiv:1309.3111*.
- Bach, F. (2017). Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53.
- Banerjee, A., Chen, T., and Zhou, Y. (2020). De-randomized pac-bayes margin bounds: Applications to non-convex and non-smooth predictors. *arXiv preprint arXiv:2002.09956*.
- Barron, A. R. and Klusowski, J. M. (2019). Complexity, statistical risk, and metric entropy of deep nets using total path variation. *arXiv preprint arXiv:1902.00800*.
- Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. (2019). Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20(63):1–17.
- Bartlett, P. L. and Mendelson, S. (2001). Rademacher and gaussian complexities: Risk bounds and structural results. In Helmbold, D. and Williamson, B., editors, *Computational Learning Theory*, volume 2111, pages 224–240. Springer Berlin Heidelberg. Series Title: Lecture Notes in Computer Science.
- Bordelon, B., Chaudhry, H. T., and Pehlevan, C. (2024). Infinite limits of multi-head transformer dynamics. *arXiv preprint arXiv:2405.15712*.
- Burer, S. and Monteiro, R. D. (2003). A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Mathematical Programming*, 95(2):329–357.
- Candès, E. J. and Plan, Y. (2011). Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Transactions on Information Theory*, 57(4):2342–2359.
- Candès, E. J. and Recht, B. (2009). Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772.

- Chandrasekaran, V., Recht, B., Parrilo, P. A., and Willsky, A. S. (2012). The convex geometry of linear inverse problems. *Foundations of Computational Mathematics*, 12(6):805–849.
- Deora, P., Ghaderi, R., Taheri, H., and Thrampoulidis, C. (2024). On the optimization and generalization of multi-head attention. *Transactions on Machine Learning Research*.
- Dziugaite, G. K. and Roy, D. M. (2017). Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *Conference on Uncertainty in Artificial Intelligence*.
- Edelman, B. L., Goel, S., Kakade, S., and Zhang, C. (2022). Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages 5793–5831. PMLR.
- Feldman, V. and Vondrak, J. (2019). High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Conference on Learning Theory*, pages 1270–1279. PMLR.
- Ge, R., Jin, C., and Zheng, Y. (2017). No spurious local minima in nonconvex low rank problems: a unified geometric analysis. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, pages 1233–1242. JMLR.org.
- Giampouras, P., Vidal, R., Rontogiannis, A., and Haeffele, B. D. (2020). A novel variational form of the Schatten- $p$  quasi-norm. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS ’20*, Red Hook, NY, USA. Curran Associates Inc.
- Golowich, N., Rakhlin, A., and Shamir, O. (2018). Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, pages 297–299. PMLR.
- Gunasekar, S., Lee, J., Soudry, D., and Srebro, N. (2018a). Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, pages 1832–1841. PMLR.
- Gunasekar, S., Lee, J. D., Soudry, D., and Srebro, N. (2018b). Implicit bias of gradient descent on linear convolutional networks. *Advances in neural information processing systems*, 31.
- Gunasekar, S., Woodworth, B. E., Bhojanapalli, S., Neyshabur, B., and Srebro, N. (2017). Implicit regularization in matrix factorization. *Advances in neural information processing systems*, 30.
- Haeffele, B. D. and Vidal, R. (2015). Global optimality in tensor factorization, deep learning, and beyond. *arXiv preprint arXiv:1506.07540*.
- Haeffele, B. D. and Vidal, R. (2017). Global optimality in neural network training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7331–7339.
- Haeffele, B. D. and Vidal, R. (2020). Structured low-rank matrix factorization: Global optimality, algorithms, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(6):1468–1482.
- HaoChen, J. Z., Wei, C., Lee, J., and Ma, T. (2021). Shape matters: Understanding the implicit bias of the noise covariance. In *Conference on Learning Theory*, pages 2315–2357. PMLR.
- Hendrickx, J. and Olshevsky, A. (2010). Matrix  $p$ -norms are np-hard to approximate if  $p \neq 1, 2, \infty$ . *SIAM Journal on Matrix Analysis and Applications*, 31(5):2802.
- Huang, C. (2020). ReLU networks are universal approximators via piecewise linear or constant functions. *Neural Computation*, 32(11):2249–2278.
- Imaizumi, M. and Schmidt-Hieber, J. (2023). On generalization bounds for deep networks based on loss surface implicit regularization. *IEEE Transactions on Information Theory*, 69(2):1203–1223.
- Jia, X., Wang, H., Peng, J., Feng, X., and Meng, D. (2023). Preconditioning matters: Fast global convergence of non-convex matrix factorization via scaled gradient descent. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S., editors, *Advances in Neural Information Processing Systems*, volume 36, pages 76202–76213. Curran Associates, Inc.
- Jin, J., Li, Z., Lyu, K., Du, S. S., and Lee, J. D. (2023). Understanding incremental learning of gradient descent: A fine-grained analysis of matrix sensing. In *Proceedings of the 40th International Conference on Machine Learning*, pages 15200–15238. PMLR. ISSN: 2640-3498.
- Kakade, S. M., Sridharan, K., and Tewari, A. (2008). On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc.
- Koltchinskii, V., Lounici, K., and Tsybakov, A. B. (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *The Annals of Statistics*, 39(5):2302–2329. Publisher: Institute of Mathematical Statistics.
- Li, G. and Wei, Y. (2023). A non-asymptotic framework for approximate message passing in spiked models. *arXiv preprint arXiv:2208.03313*.

- Li, Y., Ildiz, M. E., Papailiopoulos, D., and Oymak, S. (2023). Transformers as algorithms: Generalization and stability in in-context learning. In *Proceedings of the 40th International Conference on Machine Learning*, pages 19565–19594. PMLR.
- Li, Z., Luo, Y., and Lyu, K. (2020). Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. In *International Conference on Learning Representations*.
- Lugosi, G. and Neu, G. (2022). Generalization bounds via convex analysis. In *Conference on Learning Theory*, pages 3524–3546. PMLR.
- Ma, C., Wang, K., Chi, Y., and Chen, Y. (2020). Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution. *Foundations of Computational Mathematics*, 20(3):451–632.
- McAllester, D. A. (1999). PAC-bayesian model averaging. In *Proceedings of the twelfth annual conference on Computational learning theory*, pages 164–170. ACM.
- Muthukumar, R. and Sulam, J. (2023). Sparsity-aware generalization theory for deep neural networks. In *Proceedings of Thirty Sixth Conference on Learning Theory*, pages 5311–5342. PMLR. ISSN: 2640-3498.
- Negahban, S. and Wainwright, M. J. (2011). Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, 39(2):1069 – 1097.
- Neyshabur, B., Bhojanapalli, S., Mcallester, D., and Srebro, N. (2017). Exploring generalization in deep learning. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Neyshabur, B., Bhojanapalli, S., and Srebro, N. (2018). A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. In *International Conference on Learning Representations*.
- Nichani, E., Damian, A., and Lee, J. D. (2024). How transformers learn causal structure with gradient descent. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F., editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 38018–38070. PMLR.
- Oymak, S. and Soltanolkotabi, M. (2019). Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *Proceedings of the 36th International Conference on Machine Learning*, pages 4951–4960. PMLR. ISSN: 2640-3498.
- Recht, B., Xu, W., and Hassibi, B. (2008). Necessary and sufficient conditions for success of the nuclear norm heuristic for rank minimization. In *2008 47th IEEE Conference on Decision and Control*. IEEE.
- Reddy, T. U. K. and Vidyasagar, M. (2023). Convergence of momentum-based heavy ball method with batch updating and/or approximate gradients. In *2023 Ninth Indian Control Conference (ICC)*, pages 182–187.
- Rockafellar, R. T. (1970). *Convex Analysis*. Princeton University Press, Princeton.
- Shalev-Shwartz, S., Shamir, O., Srebro, N., and Sridharan, K. (2009). Stochastic convex optimization. In *COLT*, volume 2, page 5.
- Singh, S. S. (2023). Analyzing transformer dynamics as movement through embedding space. *arXiv preprint arXiv:2308.10874*.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. (2018). The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57.
- Stöger, D. and Soltanolkotabi, M. (2021). Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. *Advances in Neural Information Processing Systems*, 34:23831–23843.
- Stöger, D. and Zhu, Y. (2024). Non-convex matrix sensing: Breaking the quadratic rank barrier in the sample complexity. *arXiv preprint arXiv:2408.13276*.
- Team, G. (2024). Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Tian, Y., Wang, Y., Chen, B., and Du, S. S. (2023). Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. *Advances in Neural Information Processing Systems*, 36:71911–71947.
- Trauger, J. and Tewari, A. (2024). Sequence length independent norm-based generalization bounds for transformers. In *International Conference on Artificial Intelligence and Statistics*, pages 1405–1413. PMLR.
- Vapnik, V. N. (2000). *The Nature of Statistical Learning Theory*. Springer.
- Vapnik, V. N. and Chervonenkis, A. Y. (1968). The uniform convergence of frequencies of the appearance of events to their probabilities. *Doklady Akademii Nauk SSSR*, 181(4):781–783.
- Vardi, G. (2023). On the implicit bias in deep-learning algorithms. *Commun. ACM*, 66(6):86–93.

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Vershynin, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Vidal, R., Zhu, Z., and Haeffele, B. D. (2022). *Optimization Landscape of Neural Networks*, page 200–228. Cambridge University Press.
- Wen, K., Li, Z., and Ma, T. (2023). Sharpness minimization algorithms do not only minimize sharpness to achieve better generalization. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Yang, Y., Wipf, D. P., et al. (2022). Transformers from an optimization perspective. *Advances in Neural Information Processing Systems*, 35:36958–36971.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. (2021). Understanding deep learning (still) requires rethinking generalization. *Commun. ACM*, 64(3):107–115.
- Zhang, R., Frei, S., and Bartlett, P. L. (2024). Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55.
- Zhang, Z., Fang, J., Lin, J., Zhao, S., Xiao, F., and Wen, J. (2020). Improved upper bound on the complementary error function. *Electronics Letters*, 56(13):663–665.
- Zhou, Y., Alon, U., Chen, X., Wang, X., Agarwal, R., and Zhou, D. (2024). Transformers can achieve length generalization but not robustly. In *International Conference on Learning Representations*.

## Checklist

1. For all models and algorithms presented, check if you include:
  - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes/No/Not Applicable]
  - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes/No/Not Applicable]
  - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes/No/Not Applicable]
2. For any theoretical claim, check if you include:
  - (a) Statements of the full set of assumptions of all theoretical results. [Yes/No/Not Applicable]
  - (b) Complete proofs of all theoretical results. [Yes/No/Not Applicable]
  - (c) Clear explanations of any assumptions. [Yes/No/Not Applicable]
3. For all figures and tables that present empirical results, check if you include:
  - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes/No/Not Applicable]
  - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes/No/Not Applicable]
  - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes/No/Not Applicable]
  - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Yes/No/Not Applicable]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
  - (a) Citations of the creator If your work uses existing assets. [Yes/No/Not Applicable]
  - (b) The license information of the assets, if applicable. [Yes/No/Not Applicable]
  - (c) New assets either in the supplemental material or as a URL, if applicable. [Yes/No/Not Applicable]
  - (d) Information about consent from data providers/curators. [Yes/No/Not Applicable]
  - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Yes/No/Not Applicable]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
  - (a) The full text of instructions given to participants and screenshots. [Yes/No/Not Applicable]
  - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Yes/No/Not Applicable]
  - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Yes/No/Not Applicable]

---

# A Convex Relaxation Approach to Generalization Analysis for Parallel Positively Homogeneous Networks: Supplementary Materials

---

In this supplementary material, we provide a detailed discussion of the rigorous technical aspects omitted from the main text. Additionally, we present a comprehensive review of related works and include a few numerical experiments. Below is the table of contents for this appendix/supplementary material.

## Contents

<b>1</b>	<b>INTRODUCTION</b>	<b>1</b>
<b>2</b>	<b>PROBLEM FORMULATION</b>	<b>3</b>
<b>3</b>	<b>CONVEX BOUNDS FOR LEARNING</b>	<b>4</b>
<b>4</b>	<b>STATISTICAL BOUNDS</b>	<b>5</b>
<b>5</b>	<b>APPLICATIONS</b>	<b>7</b>
<b>6</b>	<b>CONCLUSIONS</b>	<b>9</b>
<b>A</b>	<b>CONVEX BOUNDS FOR LEARNING</b>	<b>14</b>
A.1	Induced Regularizer in Convex Space . . . . .	14
A.2	Proof of Theorem 1 . . . . .	15
<b>B</b>	<b>STATISTICAL BOUNDS</b>	<b>20</b>
B.1	Computing Function Class Capacities . . . . .	21
B.2	Proof of Theorem 2 . . . . .	22
<b>C</b>	<b>APPLICATIONS</b>	<b>28</b>
C.1	Low-Rank Matrix Sensing . . . . .	28
C.2	Structured Matrix Sensing . . . . .	34
C.3	Two-Layer Linear NN . . . . .	36
C.4	Two-Layer ReLU NN . . . . .	41
C.5	Multi-head Attention . . . . .	47
<b>D</b>	<b>GOODS EVENTS</b>	<b>54</b>
D.1	Concentration of Norms . . . . .	54
D.2	Concentration of Convex functions . . . . .	55
D.3	Concentration of Equilibria . . . . .	58
D.4	Concentration of Polar . . . . .	60
<b>E</b>	<b>NUMERICAL EXPERIMENTS</b>	<b>63</b>
<b>F</b>	<b>OTHER RELATED WORKS</b>	<b>63</b>

## G PRELIMINARIES

66

G.1 Convex Functions . . . . .	66
G.2 Concentration of Measure . . . . .	67

## A CONVEX BOUNDS FOR LEARNING

In this section, we discuss the proof for Theorem 1 that establishes the optimality gaps in the empirical and population landscapes. First, we analyze the convexity of the induced regularizer,  $\Omega(\cdot)$  and properties of the stationary points in non-convex landscape. These are the key components of our proof for Theorem 1. We state a more general version of Assumption 2 by having the flexibility of the loss being strongly convex to derive tighter results.

**Assumption 2'** (Convex Loss). *The loss  $\ell(Y, \hat{Y})$  is second-order differentiable (written  $\ell \in \mathcal{C}^2$ ),  $\alpha$ -strong and  $L$ -smooth w.r.t.  $\hat{Y}$ , i.e., for any  $Y, \hat{Y} \in \mathbb{R}^{n_Y}$*

$$0 \preceq \alpha I_{n_Y} \preceq \nabla_{\hat{Y}}^2 \ell(Y, \hat{Y}) \preceq L I_{n_Y}. \quad (32)$$

Additionally, the gradient of the loss is bi-Lipschitz; that is, for all  $Y_1, Y_2, \hat{Y}_1, \hat{Y}_2 \in \mathbb{R}^{n_Y}$

$$\|\nabla_{\hat{Y}} \ell(Y_2, \hat{Y}_2) - \nabla_{\hat{Y}} \ell(Y_1, \hat{Y}_1)\| \leq L [\|Y_2 - Y_1\|_2 + \|\hat{Y}_2 - \hat{Y}_1\|_2], \quad (33)$$

and the loss is constant if both the arguments are the same, i.e., for all  $Y_1, Y_2 \in \mathbb{R}^{n_Y}$ ,  $\ell(Y_1, Y_1) = \ell(Y_2, Y_2)$ .

Note that we allow  $\alpha = 0$ , in which case we recover Assumption 2.

### A.1 Induced Regularizer in Convex Space

First, we show that the induced regularizer is convex in the function spaces through Proposition 1.

**Proposition 1** (Convexity of induced regularizer). *Suppose assumptions 1-2' hold. Then  $\Omega(f)$  is convex in  $f$  in the space of functions  $\mathbb{R}^{n_X} \rightarrow \mathbb{R}^{n_Y}$ .*

*Proof.* This proof is infinite dimensional extension of Haeffele and Vidal (2015). Recall the definition of induced regularizer:

$$\Omega(f) := \inf_{r, \{W_j\}} \Theta_r(\{W_j\}) \text{ such that } f(X) = \Phi_r(\{W_j\}); \forall X \in \mathcal{X}. \quad (34)$$

Define the function class

$$\mathcal{F}_\Phi := \{\Phi_r(\{W_j\}) : r \in \mathbb{N}, W_j \in \mathcal{W}\}. \quad (35)$$

By definition if  $f \notin \mathcal{F}_\Phi$  then  $\Omega(f)$  evaluates to infinity. Now suppose that  $\beta \geq 0$  and for any  $f \in \mathcal{F}_\Phi$ ,

$$\Omega(\beta f) = \inf_{r, \{W_j\}} \Theta_r(\{W_j\}) \text{ such that } \beta f(X) = \Phi_r(\{W_j\}); \forall X \in \mathcal{X}, \quad (36)$$

Now by Assumption 2', there exists  $\hat{\beta}$  such that  $\beta \Phi_r(\{W_j\}) = \Phi_r(\{\hat{\beta} W_j\})$ , and  $\beta \Theta_r(\{W_j\}) = \Theta_r(\{\hat{\beta} W_j\})$  (throughout note that this scaling is applied only to the  $\mathcal{W}_p$  subset of parameters from Assumption 2, but we do not notate this explicitly for brevity of notation). Now we perform a change of variables in the induced regularizer, obtaining

$$\Omega(\beta f) = \inf_{r, \{\hat{\beta} W_j\}} \Theta_r(\{\hat{\beta} W_j\}) \text{ such that } \beta f(X) = \Phi_r(\{\hat{\beta} W_j\}); \forall X \in \mathcal{X}. \quad (37)$$

Then we have that

$$\Omega(\beta f) = \inf_{r, \{W_j\}} \beta \Theta_r(\{W_j\}) \text{ such that } \beta f(X) = \beta \Phi_r(\{W_j\}); \forall X \in \mathcal{X} = \beta \Omega(f). \quad (38)$$

We have established that the function  $\Omega(\cdot)$  is 1-degree homogeneous. Now we prove that the function  $\Omega(\cdot)$  is sub-additive. Choose any  $f_1, f_2 \in \mathcal{F}_\Phi$ , because the case when either of them is not in  $\mathcal{F}_\Phi$  is trivially sub-additive. Recall

$$\Omega(f_1) = \inf_{r, \{W_j\}} \Theta_r(\{W_j\}) \text{ such that } f_1(X) = \Phi_r(\{W_j\}); \forall X \in \mathcal{X}, \quad (39)$$

$$\Omega(f_2) = \inf_{r, \{W_j\}} \Theta_r(\{W_j\}) \text{ such that } f_2(X) = \Phi_r(\{W_j\}); \forall X \in \mathcal{X}, \quad (40)$$

$$\Omega(f_1 + f_2) = \inf_{r, \{W_j\}} \Theta_r(\{W_j\}) \text{ such that } f_1(X) + f_2(X) = \Phi_r(\{W_j\}); \forall X \in \mathcal{X}. \quad (41)$$

For any  $\epsilon > 0$  let  $(r_1, \{W_j^1\})$  and  $(r_2, \{W_j^2\})$  be parameters which come within  $\epsilon$  of the infimum in the optimization problems for  $\Omega(f_1)$  and  $\Omega(f_2)$  respectively. Then note that

$$\Omega(f_1 + f_2) \leq \Theta_{r_1}(\{W_j^1\}) + \Theta_{r_2}(\{W_j^2\}) \leq \Omega(f_1) + \Omega(f_2) + 2\epsilon. \quad (42)$$

Letting  $\epsilon \rightarrow 0$  gives that  $\Omega(f_1 + f_2) \leq \Omega(f_1) + \Omega(f_2)$ . Thus, as  $\Omega(\cdot)$  is both positively homogenous with degree one and sub-additive, it is convex.  $\square$

From the above proposition we have that  $\Omega(\cdot)$  is a convex function, therefore we have that  $C(\cdot)$  is indeed a convex function in the prediction functions space. Our results primarily depend upon the optimal regularization of the globally optimal solution of a convex function,  $C(\cdot)$ . As we operate in the space of functions, it is very unlikely that we have the knowledge of the global optima. Nevertheless, by exploiting the convexity of  $C(\cdot)$  we can upper bound the optimal regularization. Proposition 2 establishes the upper bound for the optimal regularization for regression loss.

**Proposition 2.** Consider  $\ell(Y_1, Y_2) = \frac{1}{2} \|Y_1 - Y_2\|_2^2$ ,  $\{W_j\} \in \mathcal{F}_W$ . Suppose  $X \sim \mu$ ,  $\epsilon$  is random variable such that  $\mathbb{E}[\epsilon] = 0$  and independent from  $x$ . Let  $Y = \Phi_r(\{W_j\})(X) + \epsilon$ , and suppose  $f_\mu^*$  is the global optimal solution of  $C_\mu(\cdot)$ . Then we have

$$\Omega(\Phi_r(\{W_j\})) \geq \Omega(f_\mu^*). \quad (43)$$

*Proof.* As  $f_\mu^*$  is the global optimal solution, we have that

$$\mathbb{E} \left[ \frac{1}{2} \|\Phi_r(\{W_j\})(X) + \epsilon - f_\mu^*(X)\|_2^2 \right] + \lambda \Omega(f_\mu) \leq \mathbb{E} \left[ \frac{1}{2} \|\Phi_r(\{W_j\})(X) + \epsilon - \Phi_r(\{W_j\})(X)\|_2^2 \right] + \lambda \Omega(\Phi_r(\{W_j\})) \quad (44)$$

Now, by re-arranging the terms we obtain

$$\mathbb{E} \left[ \frac{1}{2} \langle \Phi_r(\{W_j\})(X) - f_\mu^*(X), \epsilon \rangle \right] + \lambda \Omega(f_\mu) \leq \lambda \Omega(\Phi_r(\{W_j\})). \quad (45)$$

As  $\epsilon$  is independent of  $X$ ,

$$\frac{1}{2} \langle \mathbb{E}_X [\Phi_r(\{W_j\})(X) - f_\mu^*(X)], \mathbb{E}_\epsilon [\epsilon] \rangle + \lambda \Omega(f_\mu) \leq \lambda \Omega(\Phi_r(\{W_j\})). \quad (46)$$

Then we have

$$\Omega(f_\mu) \leq \Omega(\Phi_r(\{W_j\})). \quad (47)$$

$\square$

## A.2 Proof of Theorem 1

Optimization algorithms used to optimize DNNs try to find the set of parameters that are first-order optimal. However, we do not have a guarantee that these points are saddle/local minima/global minima. In proposition 3, we provide properties that any first-order optimal satisfies for positively homogeneous networks.

**Proposition 3** (Stationary Points). Under assumption 2, if  $\{W_j\}$  are stationary points of  $\text{NC}_\mu(\cdot)$ , then for all  $j \in [r]$ ,

$$\langle -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W_j) \rangle = \theta(W_j). \quad (48)$$

*Proof.* This proof is similar to that of Proposition 2 in Haeffele and Vidal (2020) but applied to a general class of (*slightly*) positively homogeneous functions (see assumption 2).

From assumption 2, there exists a subset of parameters where both  $\theta$  and  $\phi$  are positively homogeneous. Let  $\mathbf{w}_i$  be the subset of parameters in  $\mathcal{W}_p$  from assumption 2. Then we have

$$\langle \mathbf{w}_i, \partial_{\mathbf{w}_i} \theta(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n) \rangle = \lim_{\epsilon \rightarrow 0} \left[ \frac{\theta(\mathbf{w}_1, \dots, (1+\epsilon)\mathbf{w}_i, \dots, \mathbf{w}_n)}{\epsilon} - \frac{\theta(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n)}{\epsilon} \right]. \quad (49)$$

Let  $p_i$  be the homogeneous degree of the parameters  $\mathbf{w}_i$ . Note that  $\partial_{\mathbf{w}_i} \theta(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n) \in \mathbb{R}^{\dim(\mathbf{w}_i) \times 1}$ ,  $\partial_{\mathbf{w}_i} \phi(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n) \in \mathbb{R}^{\dim(\mathbf{w}_i) \times n_Y}$ . Then

$$\langle \mathbf{w}_i, \partial_{\mathbf{w}_i} \theta(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n) \rangle = \theta(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n) \lim_{\epsilon \rightarrow 0} \frac{(1+\epsilon)^{p_i} - 1}{\epsilon}, \quad (50)$$

$$= p_i \theta(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n). \quad (51)$$

Similarly, following a similar argument for  $\phi$  we obtain

$$\langle \partial_{\mathbf{w}_i} \phi(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n), \mathbf{w}_i \rangle = p_i \phi(\mathbf{w}_1, \dots, \mathbf{w}_i, \dots, \mathbf{w}_n). \quad (52)$$

As  $W_j$  are the stationary points we have that

$$0 \in \partial_{W_j} \Phi_r(\{W_j\}) \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\}))_\mu + \lambda \partial_{W_j} \Theta_r(\{W_j\}). \quad (53)$$

Since  $\Phi_r(\{W_j\}) = \sum_{j=1}^r \phi(W_j)$ , we have that  $\partial_{W_j} \Phi_r(\{W_j\}) = \partial_{W_j} \phi(W_j)$ . Similarly,  $\partial_{W_j} \Theta_r(\{W_j\}) = \partial_{W_j} \theta(W_j)$  holds true. Consequently,

$$0 \in \partial_{W_j} \phi(W_j) \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\}))_\mu + \lambda \partial_{W_j} \theta(W_j). \quad (54)$$

Letting  $W_j = [\mathbf{w}_1 \ \dots \ \mathbf{w}_n]$ , for all  $\mathbf{w}_i$  it holds that

$$0 \in \partial_{\mathbf{w}_i} \phi(W_j) \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\}))_\mu + \lambda \partial_{\mathbf{w}_i} \theta(W_j). \quad (55)$$

Taking the inner product of the above equation with  $\mathbf{w}_i$ , when  $p_i \neq 0$  we have that

$$0 \in \mathbf{w}_i^T \partial_{\mathbf{w}_i} \phi(W_j) \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\}))_\mu + \lambda \mathbf{w}_i^T \partial_{\mathbf{w}_i} \theta(W_j). \quad (56)$$

From (51) we have that

$$0 = p_i \phi(W_j)^T \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\}))_\mu + \lambda p_i \theta(W_j). \quad (57)$$

Rearranging, we obtain

$$\langle -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W_j) \rangle = \theta(W_j), \quad (58)$$

which holds for all  $j \in [r]$ .  $\square$

Proposition 3 establishes that at any stationary point, the inner product between the prediction errors and the predictions equates to the the current regularization. Next, we exploit this property of stationary points that enable us to tie the non-convex landscape to its convex counterpart. Lemma 1 establishes the difference between the non-convex and convex objective values at stationary points.

**Lemma 1** (Optimality Gap). *Let  $\ell(\cdot, \cdot)$  denote any  $L$ -smooth, and  $\alpha$ -strongly convex loss function, let  $q$  be some measure, and suppose that  $\{W_j\}$  is a stationary point of  $\text{NC}_q(\{W_j\})$ . Let  $f_q^*$  denote the global minimizer of  $\text{C}_q(\cdot)$ . Then for any  $f \in L^2(q)$ , we have that*

$$\text{C}_q(f_q^*) \leq \text{NC}_q(\{W_j\}) \leq \text{C}_q(f) + \lambda \Omega_q(f) \left[ \Omega_q^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2 \quad (59)$$



*Proof.* The loss  $\ell(Y, \hat{Y})$  is  $(L, \lambda)$ -convex in  $\hat{Y}$ . Therefore, for any functions  $g : \mathcal{X} \times E \rightarrow \mathcal{Y}$ , and  $f_1, f_2 \in L^2(q)$  we have that

$$\ell(g(X, \epsilon), f(X)) \geq \ell(g(X, \epsilon), \Phi_r(\{W_j\})(X)) \quad (60)$$

$$+ \langle \nabla_{\hat{Y}} \ell(g(X, \epsilon), \Phi_r(\{W_j\})(X)), f(X) - \Phi_r(\{W_j\})(X) \rangle_{\mathcal{Y}} \quad (61)$$

$$+ \frac{\alpha}{2} \|f(X) - \Phi_r(\{W_j\})(X)\|_{\mathcal{Y}}^2. \quad (62)$$

Taking expectations of both sides with respect to the probability measure  $q$ , we have that

$$\ell(g, f)_q \geq \ell(g, \Phi_r(\{W_j\}))_q + \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f - \Phi_r(\{W_j\}) \rangle_q + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2. \quad (63)$$

As the  $\{W_j\}$  are the stationary points of  $\text{NC}_q(\{W_j\})$  from Proposition 3 we have that for all  $j \in [r]$

$$\langle -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W_j) \rangle_q = \theta(W_j). \quad (64)$$

Summing the above identity up overall  $j$ , it holds that

$$\langle -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_q = \Theta_r(\{W_j\}). \quad (65)$$

Therefore, plugging this identity into the inequality (63), we have that

$$\ell(g, f)_q \geq \underbrace{\ell(g, \Phi_r(\{W_j\}))_q}_{\text{NC}_q(\{W_j\})} + \lambda \Theta_r(\{W_j\}) + \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f \rangle_q + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2, \quad (66)$$

which implies that

$$\ell(g, f)_q + \lambda \langle -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f \rangle_q \geq \text{NC}_q(\{W_j\}) + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2. \quad (67)$$

We have established from Proposition 1 that  $\Omega$  is a convex function. As a well-known result from convex analysis (see Proposition 5) we have that for any convex function  $\Omega$  and any  $f, g \in L^2(q)$ , it holds that  $\langle f, g \rangle_q \leq \Omega_q(f) \Omega_q^\circ(g)$ . Consequently,

$$\ell(g, f)_q + \lambda \Omega_q(f) \Omega_q^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) \geq \text{NC}_q(\{W_j\}) + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2. \quad (68)$$

Therefore, rearranging,

$$\underbrace{\ell(g, f)_q + \lambda \Omega_q(f)}_{\text{C}_q(f)} + \lambda \Omega_q(f) \left[ \Omega_q^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \geq \text{NC}_q(\{W_j\}) + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2, \quad (69)$$

and, as a result,

$$\text{NC}_q(\{W_j\}) \leq \text{C}_q(f) + \lambda \Omega_q(f) \left[ \Omega_q^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2. \quad (70)$$

Let  $f_q^* = \arg \min_f \text{C}_q(f)$ , and  $(r^*, \{W_j^*\}) = \arg \min_{r, \{W_j\}} \text{NC}_q(\{W_j\})$ . Since  $f_q^*$  is the minimizer of  $\text{C}_q(f)$ , it holds that

$$\text{C}_q(f_q^*) \leq \text{C}_q(\Phi_{r^*}(\{W_j^*\})) = \ell(g, \Phi_{r^*}(\{W_j^*\}))_q + \lambda \Omega_q(\Phi_{r^*}(\{W_j^*\})). \quad (71)$$

Therefore, we obtain

$$\text{C}_q(f_q^*) \leq \ell(g, \Phi_{r^*}(\{W_j^*\}))_q + \lambda \Omega_q(\Phi_{r^*}(\{W_j^*\})) \quad (72)$$

$$\leq \ell(g, \Phi_{r^*}(\{W_j^*\}))_q + \lambda \Theta_{r^*}(\{W_j^*\}) \quad (73)$$

$$= \text{NC}_q(\{W_j^*\}) \quad (74)$$

$$\leq \text{NC}_q(\{W_j\}). \quad (75)$$

Therefore, combining Equations (75) and (70), we obtain the bound

$$\mathbf{C}_q(f_q^*) \leq \mathbf{NC}_q(\{W_j\}) \leq \mathbf{C}_q(f) + \lambda \Omega_q(f) \left[ \Omega_q^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_q^2. \quad (76)$$

□

Lemma 1 has established that the non-convex objective,  $\mathbf{NC}_q(\cdot)$  is both upper and lower bounded by the convex function,  $\mathbf{C}_q(\cdot)$ . Now, we utilize this result to compute the empirical gap with the measure,  $\mu_N$  for the stationary points obtained from the ERM. On these stationary points, we bound the optimality gap by changing the measure to  $\mu$ , i.e., the behavior of ERM's first-order points on population landscape.

**Theorem 3** (Global Optimality). *Under assumptions 1, 2', 3. Let  $f_{\mu_N}^*$  (or  $f_\mu^*$ ) be the global minimizer for  $\mathbf{C}_{\mu_N}(\cdot)$  (or  $\mathbf{C}_\mu(\cdot)$ ). For any stationary points,  $(r, \{W_j\})$  of the function  $\mathbf{NC}_{\mu_N}(\cdot)$  and any  $f \in L^2(\mu) \cap L^2(\mu_N)$  the following items are true:*

1. *Empirical optimality gap:*

$$\mathbf{C}_{\mu_N}(f_{\mu_N}^*) \leq \mathbf{NC}_{\mu_N}(\{W_j\}) \leq \mathbf{C}_{\mu_N}(f) + \lambda \Omega(f) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_{\mu_N}^2, \quad (77)$$

2. *Population optimality gap:*

$$\begin{aligned} \mathbf{C}_\mu(f_\mu^*) \leq \mathbf{NC}_\mu(\{W_j\}) \leq \mathbf{C}_\mu(f) + \lambda \Omega(f) \left[ \Omega_\mu^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2 \\ + [\langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N}]. \end{aligned} \quad (78)$$

where  $\Omega_q^\circ(\cdot)$  is the polar in the measure  $q$  defined as

$$\Omega_q^\circ(g) := \sup_{\theta(W) \leq 1} \langle g, \phi(W) \rangle_q \quad (79)$$

**Remarks:** Setting  $f = f_{\mu_N}^*$  in Equation 77 and taking  $(r, \{W_j\})$  to be any stationary point of  $\mathbf{NC}_{\mu_N}(\cdot)$  gives a means to verify if  $\{W_j\}$  is a globally optimal solution. We see that it suffices to check if  $\Phi_r(\{W_j\})$  is a first-order stationary point of  $\mathbf{C}_{\mu_N}(\cdot)$ , which is a necessary condition for a local minimum of convex functions.

From convex analysis, if a function  $f \in L^2(\mu_N)$  is a first-order solution of  $\mathbf{C}_{\mu_N}$  then we have that 0 belongs to the sub-gradient of  $\mathbf{C}_\mu(\cdot)$  at  $f$ . As the loss  $\ell$  is first-order differentiable (by Assumption 3 or 2') we have that

$$0 \in \partial \mathbf{C}_{\mu_N}(f) \iff -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, f)_{\mu_N} \in \partial \Omega(f), \quad (80)$$

where  $\partial \mathbf{C}_{\mu_N}(f)$  denotes the subgradient of  $\mathbf{C}$  (viewed as a function of  $f$ ). The above condition for  $f$  can also be verified by a dual notion known as the polar condition, Definition 6 (Rockafellar, 1970). The sub-gradient of a convex function can be defined through the notion of it's polar via

$$\partial \Omega_{\mu_N}(f) = \{g \in L_2(\mu_N) : \langle g, f \rangle_{\mu_N} = \Omega_{\mu_N}(f), \Omega_{\mu_N}^\circ(g) \leq 1\}. \quad (81)$$

From Lemma 1 in the supplement of Haeffele and Vidal (2017) the following statements are equivalent:

1.  $\{W_j\}$  is an optimal factorization of  $f$ ; i.e,  $\Theta_r(\{W_j\}) = \Omega_{\mu_N}(f)$ .
2.  $\exists h \in L^2(\mu_N)$  such that  $\Omega_{\mu_N}^\circ(h) \leq 1$  and  $\langle h, \Phi_r(\{W_j\}) \rangle_{\mu_N} = \Theta_r(\{W_j\})$ .
3.  $\exists h \in L^2(\mu_N)$  such that  $\Omega_{\mu_N}^\circ(h) \leq 1$  and  $\langle h, \phi(W_j) \rangle_{\mu_N} = \theta(W_j); \forall i \in [r]$ .

Further, if (2) or (3) above is satisfied then we have that  $h \in \partial \Omega_{\mu_N}(f)$ . From Proposition 3 we have that for any stationary point  $(r, \{W_j\})$  of  $\mathbf{NC}_{\mu_N}$ ,

$$\langle -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell_{\mu_N}(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N} = \Theta_r(\{W_j\}). \quad (82)$$

Consequently, to check if a stationary point is globally optimal, it then suffices to check whether the polar condition  $\Omega_{\mu_N}^\circ(-\frac{1}{\lambda}\nabla\ell(g, \Phi_r(\{W_j\}))_{\mu_N}) \leq 1$  holds at the stationary point,  $(r, \{W_j\})$ . In the case when the polar condition holds true, the upper bound evaluates to  $C_{\mu_N}(f_\mu^*)$  matching the lower bound of  $\text{NC}_{\mu_N}(\cdot)$ , which in turn implies global optimality.

Then, we can claim the following:

At a stationary point  $\{W_j\}$ , if  $\Omega_{\mu_N}^\circ(-\frac{1}{\lambda}\nabla\ell(g, \Phi_r(\{W_j\}))_{\mu_N}) \leq 1$ , then  $\{W_j\}$  is globally optimal.

Now we prove Theorem 1.

*Proof.* The proof sketch is similar to Proposition 4 from Haeffele and Vidal (2020). Equation (12) can be obtained from the Lemma 1, for any stationary points,  $(r, \{W_j\})$  of  $\text{NC}_{\mu_N}(\cdot)$ .

Since,  $f \in L^2(\mu) \cap L^2(\mu_N) \subseteq L^2(\mu_n)$ , and the parameters satisfy the equality in Lemma 1, we can conclude that Equation (12) holds. The local minima of  $\text{NC}_{\mu_N}(\cdot)$  need not be local minima of  $\text{NC}_\mu(\cdot)$ , therefore we shall obtain an discrepancy term. From the fact that  $\ell$  is a  $\alpha$ -strongly convex function we have the inequality

$$\ell(g, f)_\mu \geq \ell_\mu(g, \Phi_r(\{W_j\})) + \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f - \Phi_r(\{W_j\}) \rangle_\mu + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2. \quad (83)$$

Adding  $\lambda\Theta_r(\{W_j\})$  on both sides we obtain the inequality

$$\ell(g, f)_\mu + \lambda\Theta_r(\{W_j\}) \geq \ell_\mu(g, \Phi_r(\{W_j\})) + \lambda\Theta_r(\{W_j\}) \quad (84)$$

$$+ \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f - \Phi_r(\{W_j\}) \rangle_\mu + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2. \quad (85)$$

Now replacing the first term on the side with  $\text{NC}_\mu(\{W_j\})$  we obtain

$$\ell(g, f)_\mu + \lambda\Theta_r(\{W_j\}) \geq \text{NC}_\mu(\{W_j\}) + \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f - \Phi_r(\{W_j\}) \rangle_\mu + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2. \quad (86)$$

From Proposition 3 we have that for stationary points  $\{W_j\}$ , it holds that  $\Theta_r(\{W_j\}) = \langle -\frac{1}{\lambda}\nabla_{\hat{Y}} \ell_{\mu_N}(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N}$ . Therefore, by plugging this into the inequality above, we obtain that

$$\ell(g, f)_\mu + \lambda \langle -\frac{1}{\lambda}\nabla_{\hat{Y}} \ell_{\mu_N}(g, \Phi_r(\{W_j\}))_{\mu_N}, \Phi_r(\{W_j\}) \rangle \quad (87)$$

$$\geq \text{NC}_\mu(\{W_j\}) + \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f - \Phi_r(\{W_j\}) \rangle_\mu + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2. \quad (88)$$

Rearranging the terms we have

$$\ell(g, f)_\mu + \lambda \langle -\frac{1}{\lambda}\nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), f \rangle_\mu \quad (89)$$

$$+ \langle \nabla_{\hat{Y}} \ell_\mu(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell_{\mu_N}(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N} \quad (90)$$

$$\geq \text{NC}_\mu(\{W_j\}) + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2. \quad (91)$$

Next, the following inequality always holds:

$$\langle f, -\frac{1}{\lambda}\nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \rangle_\mu \leq \Omega_\mu(f) \Omega_\mu^\circ \left( -\frac{1}{\lambda}\nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right). \quad (92)$$

Rearranging (98) and plugging in (92), we obtain the inequality

$$\ell(g, f)_\mu + \lambda \Omega_\mu(f) \Omega_\mu^\circ \left( -\frac{1}{\lambda}\nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) \quad (93)$$

$$+ \langle \nabla_{\hat{Y}} \ell_\mu(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell_{\mu_N}(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N} \quad (94)$$

$$\geq \text{NC}_\mu(\{W_j\}) + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2. \quad (95)$$

We add and subtract  $\Omega(f)$  to obtain

$$\ell(g, f)_\mu + \lambda \Omega_\mu(f) + \lambda \Omega_\mu(f) \left[ \Omega_\mu^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \quad (96)$$

$$+ \langle \nabla_{\hat{Y}} \ell_\mu(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell_{\mu_N}(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N} \quad (97)$$

$$\geq \text{NC}_\mu(\{W_j\}) + \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2. \quad (98)$$

Rearranging the right most term and using the definition of  $C_\mu(f)$  we obtain,

$$\text{NC}_\mu(\{W_j\}) \leq C_\mu(f) + \lambda \Omega_\mu(f) \left[ \Omega_\mu^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2 \quad (99)$$

$$+ [\langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu]. \quad (100)$$

This yields the right hand side of (13). As for the left hand side, by definition we have that for any  $(r, \{W_j\})$ ,  $C_\mu(f_\mu^*) \leq \text{NC}_\mu(\{W_j\})$ . This completes the proof.  $\square$

Theorem 1 provides the behavior of ERM solutions in the population landscape. This paves a path to bound the empirical and population objectives at these stationary points.

## B STATISTICAL BOUNDS

This section provides a more general version of Theorem 2 that does not need Assumption 7 to hold uniformly for all the data points,  $(X, Y)$ . Rather, we relax the assumption to the following.

**Assumption 7'** (Probabilistic boundedness). *There exists a convex set,  $\mathcal{C} \subseteq \mathbb{R}^{n_X} \times \mathbb{R}^{n_Y}$  such that*

$$P(\cap_{i=1}^N (X_i, \epsilon_i) \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}. \quad (101)$$

For all  $(X, \epsilon) \in \mathcal{C}$ , and  $\{W_j\} \in \mathcal{F}_{\mathcal{W}}$  the predictions and gradients are bounded; i.e.,

$$\|\Phi_r(\{W_j\})(X)\| \leq B_\Phi, \quad \|\nabla_{\hat{Y}} \ell(g(X, \epsilon), \Phi_r(\{W_j\}))(X)\| \leq B_\ell. \quad (102)$$

Further, for any  $\{W_j\}, \{\tilde{W}_j'\} \in \mathcal{F}_{\mathcal{W}}$ ,  $W, \tilde{W} \in \mathcal{F}_\theta$ ,  $(X, \epsilon) \in \mathcal{C}$ , the network  $\phi$  and  $\Phi_r$  are Lipschitz in the parameters; i.e.,

$$\|\Phi_r(\{W_j\})(X) - \Phi_r(\{\tilde{W}_j'\})(X)\| \leq \tilde{L}_\Phi \max_j \|W_j - \tilde{W}_j'\|_2, \quad (103)$$

and

$$\|\phi(W)(X) - \phi(\tilde{W})(X)\| \leq \tilde{L}_\phi \|W - \tilde{W}\|_2. \quad (104)$$

Additionally, define the quantity

$$B(\mathcal{C}) := \left[ (1 + \alpha) \sup_{\{W_j\} \in \mathcal{F}_{\mathcal{W}}} \left| \|f_\mu^* \circ \mathcal{P}_{\mathcal{C}} - \Phi_r(\{W_j\}) \circ \mathcal{P}_{\mathcal{C}}\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \right] \quad (105)$$

$$+ \sup_{\{W_j\} \in \mathcal{F}_{\mathcal{W}}, W' \in \mathcal{F}_\theta} |\langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_{\mathcal{C}}, \Phi_r(\{W_j\}) \circ \mathcal{P}_{\mathcal{C}}), \phi(W') \circ \mathcal{P}_{\mathcal{C}} \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu | \quad (106)$$

$$+ \sup_{\{W_j\} \in \mathcal{F}_{\mathcal{W}}} |\langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_{\mathcal{C}}, \Phi_r(\{W_j\}) \circ \mathcal{P}_{\mathcal{C}}), \Phi_r(\{W_j\}) \circ \mathcal{P}_{\mathcal{C}} \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu | \quad (107)$$

where  $\mathcal{P}_{\mathcal{C}}(\cdot)$  is the Euclidean projection to the set  $\mathcal{C}$ .

**Comparison with Assumption 7:** Unlike in Assumption 7, we do not require the equations (102), (103), and (104) to hold for all the inputs. However, we relax this restriction by assuming that there exists a convex set,  $\mathcal{C}$  which consists of the data points with probability at least  $1 - \delta_{\mathcal{C}}$ . For well-behaved probability distributions like sub-Gaussian distributions (see Assumption 2'), such a convex exists with very high probability; i.e., very small  $\delta_{\mathcal{C}}$ .

Now we state the general master theorem that relies on Assumptions 1, 2', 3, 4, 5, 6 and 7' (but not on Assumption 7).

**Theorem 4** (General Master Theorem). *Suppose Assumptions 1, 2', 3, 4, 5, 6 and 7' hold. Let  $\delta \in (0, 1]$  be fixed, and let  $f_\mu^*$  be the global optimum of  $\mathcal{C}_\mu$ . Suppose that  $\gamma \geq \Omega(f_\mu^*)L_\phi$ , and define*

$$\epsilon_1 = 16\gamma^2\sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{Lip}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}; \quad (108)$$

$$\epsilon_2 = 4\tilde{L}_\Phi B_\Phi \max \left\{ 1, 2L + \frac{2B_\ell}{B_\Phi}, 8\Omega(f_\mu^*) \frac{B_\ell \tilde{L}_\phi}{\tilde{L}_\Phi B_\Phi}, 8L\Omega(f_\mu^*) \right\}. \quad (109)$$

Let  $\{W_j\}$  denote any stationary point of  $\text{NC}_{\mu_N}(\cdot)$ . Then with probability at least  $1 - (\delta + \delta_C)$ , it holds that

$$\frac{1}{n_Y} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \quad (110)$$

$$\lesssim \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2n_Y} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (111)$$

$$+ \frac{B(\mathcal{C})}{n_Y} + (1 + \alpha)\epsilon_1 \sqrt{\frac{R \cdot \dim(\mathcal{W}) \log(\gamma\epsilon_2 r_\theta / L_\phi) \log(N) + \log(1/\delta)}{N}}. \quad (112)$$

**Additional Remarks:** In addition to the discussion in Section 4, the general version mentioned above (i) takes into account unbounded sub-Gaussian distributions, and (ii) imposes a weaker notion of Lipschitz continuity on the parameters. For sub-Gaussian inputs, one may choose the convex set  $\mathcal{C}$  to be a ball with radius  $\mathbb{B}(g)$ . As we grow  $g$ , the term  $B(\mathcal{C})$  decays exponentially, and  $\epsilon_2$  only grows in the order of polynomial. This fast decay allows us to keep the statistical error under control while pertaining to the optimal sample complexity. Theorem 4 reduces to Theorem 2 by setting  $\alpha = 0$  and  $\mathcal{C} = \text{conv}(\mathcal{X})$ . Under this choice, Assumption 7' coincides with Assumption 7.

We discuss the proof in section B.2. Before diving into the proof, we discuss a few preliminaries on the covering number essential to estimate the capacity of the hypotheses class.

## B.1 Computing Function Class Capacities

**Lemma 2** (Covering number of  $\mathcal{F}_\mathcal{W}$ ). *Under assumption 6, and 7' the  $\nu$ -net covering number of the set  $\mathcal{F}_\mathcal{W}$  on the metric,  $\|\cdot\|_{\infty,d}$  is upper bounded via*

$$\mathcal{C}_{\mathcal{F}_\mathcal{W}}(\nu) \leq (\mathcal{C}_\theta(L_\phi \nu / \gamma))^R, \quad (113)$$

where  $\mathcal{C}_\theta(\nu) := \mathcal{N}(\{W : \theta(W) \leq 1\}, d(\cdot, \cdot), \nu)$ .

*Proof.* Recall that

$$\mathcal{C}_{\mathcal{F}_\mathcal{W}}(\nu) := \mathcal{N}(\mathcal{F}_\mathcal{W}, \max_i d(\cdot, \cdot), \nu); \quad (114)$$

$$\mathcal{C}_\theta(\nu) := \mathcal{N}(\mathcal{F}_\theta, d(\cdot, \cdot), \nu). \quad (115)$$

By the definition of  $\nu$  covering number,

$$\mathcal{N}(\mathcal{F}_\mathcal{W}, \|\cdot\|_{\infty,d}, \nu) := \inf \left| \left\{ \{W_j^0\} \in \mathcal{F}_\mathcal{W} : \forall \{W_j\} \in \mathcal{F}_\mathcal{W} : \|\{W_j\} - \{W_j^0\}\|_{\infty,d} = \max_j d(W_j, W_j^0) \leq \nu \right\} \right|; \quad (116)$$

$$\mathcal{N}(\mathcal{F}_\theta, d(\cdot, \cdot), \nu) := \inf \left| \left\{ W^0 \in \mathcal{F}_\theta : \forall W \in \mathcal{F}_\theta : d(W, W^0) \leq \nu \right\} \right|. \quad (117)$$

Therefore we can upper bound  $\mathcal{N}(\mathcal{F}_\mathcal{W}, \|\cdot\|_{\infty,d}, \nu)$  with the product of  $\mathcal{N}(\mathcal{F}_\theta, d(\cdot, \cdot), \nu)$   $R$  times. We have

$$\mathcal{N}(\mathcal{F}_\mathcal{W}, \|\cdot\|_{\infty,d}, \nu) \leq \left[ \mathcal{N} \left( \frac{\gamma}{L_\phi} \mathcal{F}_\theta, d(\cdot, \cdot), \nu \right) \right]^R, \quad (118)$$

Rewriting the above for appropriately chosen  $\nu$  we get

$$\mathcal{N}(\mathcal{F}_\mathcal{W}, \|\cdot\|_{\infty,d}, \nu) \leq \left[ \mathcal{N} \left( \mathcal{F}_\theta, d(\cdot, \cdot), \frac{L_\phi \nu}{\gamma} \right) \right]^R. \quad (119)$$

This concludes our proof.  $\square$

**Lemma 3** (Bounding covering number). *Consider a metric space,  $(\mathcal{W} \subseteq \mathbb{R}^n, \|\cdot\|_2)$  and a compact set,  $\mathcal{F}_{\mathcal{W}} \subseteq \mathcal{W}$ . Suppose that there exist  $r < \infty$  such that  $\mathcal{F}_{\mathcal{W}} \subseteq \mathbb{B}(r)$ . Then we have*

$$\mathcal{N}(\mathcal{F}_{\mathcal{W}}, \|\cdot\|_2, \nu) \leq \left(1 + \frac{2r}{\nu}\right)^n. \quad (120)$$

*Proof.* We have that  $\mathcal{F}_{\mathcal{W}} \subseteq \mathbb{B}(r)$ . By monotonicity of covering numbers, we have that

$$\mathcal{N}(\mathcal{F}_{\mathcal{W}}, \|\cdot\|_2, \nu) \leq \mathcal{N}(\mathbb{B}(r), \|\cdot\|_2, \nu). \quad (121)$$

From Corollary 4.2.13 in Vershynin (2018) we have that,

$$\mathcal{N}(\mathcal{F}_{\mathcal{W}}, \|\cdot\|_2, \nu) \leq \mathcal{N}(\mathbb{B}(r), \|\cdot\|_2, \nu) \leq \left(1 + \frac{2r}{\nu}\right)^n. \quad (122)$$

□

## B.2 Proof of Theorem 2

This section discusses the proof of Theorem 4. We extensively use concentration results from Section D that are preliminaries for the upcoming technical details.

*Proof.* First, we recall the definition of generalization error:

$$\text{Generalization Error} := |\text{NC}_{\mu}(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})|. \quad (123)$$

We can bound the above from the optimality gaps obtained in Theorem 1 via the following decomposition:

$$\text{NC}_{\mu}(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\}) = \underbrace{[\text{NC}_{\mu}(\{W_j\}) - \text{C}_{\mu}(f_{\mu})]}_{\text{Population Gap}} - \underbrace{[\text{NC}_{\mu_N}(\{W_j\}) - \text{C}_{\mu_N}(f_{\mu_N})]}_{\text{Empirical Gap}} \quad (124)$$

$$+ \underbrace{[\text{C}_{\mu}(f_{\mu}) - \text{C}_{\mu_N}(f_{\mu_N})]}_{\text{Convex Gap}}. \quad (125)$$

From Theorem 1 we have that for any  $f_{\mu}, f_{\mu_N} \in L^2(\mu) \cap L^2(\mu_N)$  and stationary points  $(r, \{W_j\})$ , the empirical gap is bounded by

$$\text{C}_{\mu_N}(f_{\mu_N}^*) - \text{C}_{\mu_N}(f_{\mu_N}) \leq \text{NC}_{\mu_N}(\{W_j\}) - \text{C}_{\mu_N}(f_{\mu_N}) \quad (126)$$

$$\leq \lambda\Omega(f_{\mu_N}) \left[ \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f_{\mu_N} - \Phi_r(\{W_j\})\|_{\mu_N}^2, \quad (127)$$

and the population gap is bounded by

$$\text{C}_{\mu}(f_{\mu}^*) - \text{C}_{\mu}(f_{\mu}) \leq \text{NC}_{\mu}(\{W_j\}) - \text{C}_{\mu}(f_{\mu}) \quad (128)$$

$$\leq \lambda\Omega(f_{\mu}) \left[ \Omega_{\mu}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f_{\mu} - \Phi_r(\{W_j\})\|_{\mu}^2 \quad (129)$$

$$+ [\langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu} - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N}]. \quad (130)$$

For any  $f_{\mu}, f_{\mu_N}$ , subtracting the above two equations we obtain

$$\text{C}_{\mu}(f_{\mu}^*) - \text{C}_{\mu_N}(f_{\mu_N}) - \lambda\Omega(f_{\mu_N}) \left[ \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] + \frac{\alpha}{2} \|f_{\mu_N} - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (131)$$

$$\leq \text{NC}_{\mu}(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\}) \quad (132)$$

$$\leq \lambda\Omega(f_{\mu}) \left[ \Omega_{\mu}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f_{\mu} - \Phi_r(\{W_j\})\|_{\mu}^2 \quad (133)$$

$$+ [\langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu} - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N}] \quad (134)$$

$$+ [\text{C}_{\mu}(f_{\mu}) - \text{C}_{\mu_N}(f_{\mu_N}^*)]. \quad (135)$$

By choosing  $f_\mu = f_{\mu_N} = f_\mu^*$  (as  $f_\mu^* \in L^2(\mu) \cap L^2(\mu_N)$ ) and noting that  $f_\mu^*$  is not a random variable unlike  $f_{\mu_N}^*$  (which depends on the data points) we get

$$C_\mu(f_\mu^*) - C_{\mu_N}(f_\mu^*) - \lambda\Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] + \frac{\alpha}{2} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (136)$$

$$\leq \text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\}) \quad (137)$$

$$\leq \lambda\Omega(f_\mu^*) \left[ \Omega_\mu^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \quad (138)$$

$$+ [\langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N}] \quad (139)$$

$$+ [C_\mu(f_\mu^*) - C_{\mu_N}(f_{\mu_N}^*)]. \quad (140)$$

Since  $f_\mu^*$  is the global minimizer of  $C_\mu(\cdot)$ , it always holds that  $C_\mu(f_\mu^*) \leq C_\mu(f_{\mu_N}^*)$ . We use this fact to upper bound the right side term, upon which we obtain the bound

$$C_\mu(f_\mu^*) - C_{\mu_N}(f_\mu^*) - \lambda\Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] + \frac{\alpha}{2} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (141)$$

$$\leq \text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\}) \quad (142)$$

$$\leq \lambda\Omega(f_\mu^*) \left[ \Omega_\mu^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \quad (143)$$

$$+ [\langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N}] \quad (144)$$

$$+ [C_\mu(f_{\mu_N}^*) - C_{\mu_N}(f_{\mu_N}^*)]. \quad (145)$$

Now we add and subtract  $\Omega_\mu^\circ(\cdot)^5$  and  $\frac{\alpha}{2} \|f - \Phi_r(\{W_j\})\|_\mu^2$  on the right side. We then have that

$$\underbrace{C_\mu(f_\mu^*) - C_{\mu_N}(f_\mu^*)}_{=:T_1} - \lambda\Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \quad (146)$$

$$+ \frac{\alpha}{2} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (147)$$

$$\leq \text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\}) \quad (148)$$

$$\leq \lambda\Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (149)$$

$$+ \frac{\alpha}{2} \left[ \underbrace{\|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2}_{=:T_2} \right] \quad (150)$$

$$+ \lambda\Omega(f_\mu^*) \left[ \underbrace{\Omega_\mu^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right)}_{=:T_3} \right] \quad (151)$$

$$+ \left[ \underbrace{\langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_{\mu_N}}_{=:T_5} \right] \quad (152)$$

$$+ \left[ \underbrace{C_\mu(f_{\mu_N}^*) - C_{\mu_N}(f_{\mu_N}^*)}_{=:T_5} \right]. \quad (153)$$

Now we apply uniform concentration on the quantities  $T_1, T_2, T_3, T_4$ , and  $T_5$  to get bound the statistical error terms.

From assumption 7' we assume that  $\mathcal{C}$  is some convex set in  $\mathbb{R}^{n_x} \times \mathbb{R}^{n_y}$  such that the following hold true:

<sup>5</sup>We are ignoring the input arguments for brevity.

1. For any i.i.d. samples  $\{X_i, \epsilon_i\}$  the  $P(\bigcap_{i=1}^N (X_i, \epsilon_i) \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}$ .
2. For all  $(X, \epsilon) \in \mathcal{C}$  and  $\forall \zeta \in \mathcal{F}_{\mathcal{W}} : \|f_{\zeta}(X)\| \leq B_{\Phi}$ .
3. For all  $(X, \epsilon) \in \mathcal{C}$  we have  $\forall \zeta \in \mathcal{F}_{\mathcal{W}} : \|\nabla_{\hat{Y}} \ell(g(X, \epsilon), f_{\zeta}(X))\| \leq B_{\ell}$ .
4. For all  $(X, \epsilon) \in \mathcal{C}$  and  $\forall \zeta, \zeta' \in \mathcal{F}_{\mathcal{W}} : \|f_{\zeta}(X) - f_{\zeta'}(X)\| \leq \tilde{L}_{\Phi} \|\zeta - \zeta'\|_{\infty, 2}$ .
5. For all  $(X, \epsilon) \in \mathcal{C}$  and  $\forall \zeta, \zeta' \in \mathcal{F}_{\theta} : \|f_{\zeta}(X) - f_{\zeta'}(Z)\| \leq \tilde{L}_{\phi} \|\zeta - \zeta'\|_2$ .
6. For any  $\hat{Y}_1, \hat{Y}_2 \in \mathbb{R}^{n_Y}$  we have  $\|\nabla_{\hat{Y}} \ell(Y, \hat{Y}_1) - \nabla_{\hat{Y}} \ell(Y, \hat{Y}_2)\| \leq L \|\hat{Y}_1 - \hat{Y}_2\|$ .
7.  $B_{nrm}(\mathcal{C}) := \sup_{\zeta \in \mathcal{F}_{\mathcal{W}}} \left| \|f_{\mu}^* \circ \mathcal{P}_{\mathcal{C}} - f_{\zeta} \circ \mathcal{P}_{\mathcal{C}}\|_{\mu}^2 - \|f_{\mu}^* - f_{\zeta}\|_{\mu}^2 \right| < \infty$ .
8.  $B_{plr}(\mathcal{C}) := \sup_{\zeta \in \mathcal{F}_{\mathcal{W}}, \zeta' \in \mathcal{F}_{\theta}} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_{\mathcal{C}}, f_{\zeta} \circ \mathcal{P}_{\mathcal{C}}), f_{\zeta'} \circ \mathcal{P}_{\mathcal{C}} \rangle_{\mu} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu} \right| < \infty$ .
9.  $B_{eql}(\mathcal{C}) := \sup_{\zeta \in \mathcal{F}_{\mathcal{W}}} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_{\mathcal{C}}, f_{\zeta} \circ \mathcal{P}_{\mathcal{C}}), f_{\zeta} \circ \mathcal{P}_{\mathcal{C}} \rangle_{\mu} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta} \rangle_{\mu} \right| < \infty$ .

Next, we define the events

$$\begin{aligned} \mathcal{E}_{cvx}(\epsilon) &:= \{\forall \zeta \in \mathcal{F}_{\mathcal{W}} : |\mathbf{C}_{\mu_N}(f_{\zeta}) - \mathbf{C}_{\mu}(f_{\zeta})| \leq \epsilon + B_{nrm}(\mathcal{C})\}; \\ \mathcal{E}_{eql}(\epsilon) &:= \{\forall \zeta \in \mathcal{F}_{\mathcal{W}} : |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta} \rangle_{\mu} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta} \rangle_{\mu_N}| \leq \epsilon + B_{eql}(\mathcal{C})\}. \end{aligned}$$

Since  $\Omega^{\circ}(\cdot)$  is positively homogeneous function we can ignore the scalar  $-\frac{1}{\lambda}$  while defining the events below:

$$\mathcal{E}_{plr}(\epsilon) := \{\forall \zeta \in \mathcal{F}_{\mathcal{W}} : |\Omega_{\mu_N}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) - \Omega_{\mu}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta}))| \leq \epsilon + B_{plr}(\mathcal{C})\}; \quad (154)$$

$$\mathcal{E}_{nrm}(\epsilon) := \{\forall \zeta \in \mathcal{F}_{\mathcal{W}} : \left| \|f_{\mu}^* - f_{\zeta}\|_{\mu_N}^2 - \|f_{\mu}^* - f_{\zeta}\|_{\mu}^2 \right| \leq \epsilon + B_{nrm}(\mathcal{C})\}. \quad (155)$$

Finally, define the following good event:

$$\mathcal{E}_{good}(\epsilon) := \mathcal{E}_{cvx}\left(\frac{\epsilon}{4}\right) \cap \mathcal{E}_{eql}\left(\frac{\epsilon}{4}\right) \cap \mathcal{E}_{plr}\left(\frac{\epsilon}{4\lambda\Omega(f_{\mu}^*)}\right) \cap \mathcal{E}_{nrm}\left(\frac{\epsilon}{2}\right). \quad (156)$$

When the event  $\mathcal{E}_{good}(\epsilon)$  holds then we obtain following from the inequality (153),

$$-\epsilon/4 - B_{cvx}(\mathcal{C}) - \lambda\Omega(f_{\mu}^*) \left[ \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] + \frac{\alpha}{2} \|f_{\mu}^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (157)$$

$$\leq \mathbf{NC}_{\mu}(\{W_j\}) - \mathbf{NC}_{\mu_N}(\{W_j\}) \quad (158)$$

$$\leq \lambda\Omega(f_{\mu}^*) \left[ \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] - \frac{\alpha}{2} \|f_{\mu}^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 \quad (159)$$

$$+ \frac{\alpha}{2} [\epsilon/2 + B_{nrm}(\mathcal{C})] + \lambda\Omega(f_{\mu}^*) [\epsilon/(4\lambda\Omega(f_{\mu}^*)) + B_{plr}(\mathcal{C})] \quad (160)$$

$$+ [\epsilon/4 + B_{eql}(\mathcal{C})] + [\epsilon/4 + B_{nrm}(\mathcal{C})]. \quad (161)$$

For  $\alpha \geq 0$ , these inequalities imply that

$$|\mathbf{NC}_{\mu}(\{W_j\}) - \mathbf{NC}_{\mu_N}(\{W_j\})| \leq \lambda\Omega(f_{\mu}^*) \left[ \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \quad (162)$$

$$- \frac{\alpha}{2} \|f_{\mu}^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 + (1 + \alpha)\epsilon + (1 + \alpha)B_{nrm}(\mathcal{C}) + B_{eql}(\mathcal{C}) + \lambda\Omega(f_{\mu}^*)B_{plr}(\mathcal{C}). \quad (163)$$

Equation (163) holds with probability  $\mathbb{P}(\mathcal{E}_{good}(\epsilon))$ . We can bound the good event with union bound via

$$\mathbb{P}(\mathcal{E}_{good}(\epsilon)) \geq 1 - \underbrace{\mathbb{P}\left(\mathcal{E}_{nrm}^c\left(\frac{\epsilon}{2}\right)\right)}_{\text{Lemma 8}} - \underbrace{\mathbb{P}\left(\mathcal{E}_{cvx}^c\left(\frac{\epsilon}{4}\right)\right)}_{\text{Lemma 9}} - \underbrace{\mathbb{P}\left(\mathcal{E}_{eql}^c\left(\frac{\epsilon}{4}\right)\right)}_{\text{Lemma 10}} - \underbrace{\mathbb{P}\left(\mathcal{E}_{plr}^c\left(\frac{\epsilon}{4\lambda\Omega(f_{\mu}^*)}\right)\right)}_{\text{Lemma 11}} \quad (164)$$

Under Assumptions 1-6 and 7' we can apply Lemma 8, 9, 10, and 11 to bound the probability of the occurrence of the events,  $\mathcal{E}_{cvx}(\cdot)$ ,  $\mathcal{E}_{eql}(\cdot)$ ,  $\mathcal{E}_{plr}(\cdot)$ , and  $\mathcal{E}_{nrm}(\cdot)$ .



Define the constants

$$B_1 := 4n_Y L \left[ (\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right]; \quad (165)$$

$$B_2 := 16n_Y \gamma \|\nabla_{\hat{Y}} \ell\|_{\text{Lip}} \sigma_X \sqrt{(\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{E|X}^2}; \quad (166)$$

$$B_3 := 16n_Y \Omega(f_\mu^*) L_\phi \|\nabla_{\hat{Y}} \ell\|_{\text{Lip}} \sigma_X \sqrt{(\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{E|X}^2}; \quad (167)$$

$$B_4 := 128n_Y \gamma^2 \sigma_X^2; \quad (168)$$

$$\epsilon_0 := \max\{B_1, B_2, B_3, B_4\}; \quad (169)$$

$$\epsilon_1 := \max\{B_1, B_2, B_3, B_4\}; \quad (170)$$

$$b_1 := 8B_\ell \tilde{L}_\Phi; \quad (171)$$

$$b_2 := 8\tilde{L}_\Phi [B_\ell + B_\Phi L]; \quad (172)$$

$$b_3 := 32\Omega(f_\mu^*) \max\{\tilde{L}_\Phi B_\ell, L\tilde{L}_\Phi B_\Phi\}; \quad (173)$$

$$b_4 := 4\tilde{L}_\Phi B_\Phi; \quad (174)$$

$$\epsilon_2 := \max\{b_1, b_2, b_3, b_4\}. \quad (175)$$

Under the above conditions by Lemma 8, for any  $\epsilon \in [0, B_4]$  we have that

$$\mathbb{P} \left( \mathcal{E}_{nrm}^c \left( \frac{\epsilon}{2} \right) \right) \leq \delta_C + c_4 \exp \left( \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_4} \right) \right) - N \left( \frac{\epsilon}{B_4} \right)^2 \right). \quad (176)$$

By Lemma 9, for any  $\epsilon \in [0, B_1]$ , we have

$$\mathbb{P} \left( \mathcal{E}_{cvx}^c \left( \frac{\epsilon}{4} \right) \right) \leq \delta_C + 2 \exp \left( \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_1} \right) \right) - c_1 N \left( \frac{\epsilon}{B_1} \right)^2 \right), \quad (177)$$

for some positive constant,  $c_1$ .

Additionally by Lemma 10, for any  $\epsilon \in [0, B_2]$  we have that

$$\mathbb{P} \left( \mathcal{E}_{eqt}^c \left( \frac{\epsilon}{4} \right) \right) \leq \delta_C + c_1 \exp \left( \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_2} \right) \right) - N \left( \frac{\epsilon}{B_2} \right)^2 \right), \quad (178)$$

for some positive constant,  $c_2$ . Furthermore, by Lemma 11, for any  $\epsilon \in [0, B_3]$  we have that

$$\mathbb{P} \left( \mathcal{E}_{plr}^c \left( \frac{\epsilon}{4\Omega(f_\mu^*)} \right) \right) \leq \delta_C + c_3 \exp \left( \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_3} \right) \right) + \log \left( \mathcal{C}_{\mathcal{F}_\theta} \left( \frac{\epsilon}{b_3} \right) \right) - N \left( \frac{\epsilon}{B_3} \right)^2 \right), \quad (179)$$

for some positive constant,  $c_3$ .

For the inequalities (177), (178), (179), and (176) to all hold we choose  $\epsilon \in [0, \epsilon_0]$  and we upper bound the covering numbers  $\mathcal{C}_{\mathcal{F}_W}(\nu)$  as they are strictly decreasing in  $\nu$  by definition. Therefore, we have that

$$\max \left\{ \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_1} \right) \right), \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_2} \right) \right), \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_3} \right) \right), \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{b_4} \right) \right) \right\} \leq \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{\epsilon_2} \right) \right), \quad (180)$$

and

$$\log \left( \mathcal{C}_{\mathcal{F}_\theta} \left( \frac{\epsilon}{b_3} \right) \right) \leq \log \left( \mathcal{C}_{\mathcal{F}_\theta} \left( \frac{\epsilon}{\epsilon_2} \right) \right). \quad (181)$$

Now we plug in inequalities (177), (178), (179), and (176) in the inequality (164). Denote

$$c_5 := \max\{2, c_2, c_3, c_4\}. \quad (182)$$

Then

$$\mathbb{P}(\mathcal{E}_{good}(\epsilon)) \geq 1 - c_5 \exp(\log(\mathcal{C}_{\mathcal{F}_W}(\epsilon/\epsilon_3))) \times \left[ \exp\left(-c_1 N \left(\frac{\epsilon}{B_1}\right)^2\right) + \exp\left(-N \left(\frac{\epsilon}{B_2}\right)^2\right) \right] \quad (183)$$

$$+ \exp\left(\log\left(C_{\mathcal{F}_\theta}\left(\frac{\epsilon}{\epsilon_3}\right)\right) - N \left(\frac{\epsilon}{B_3}\right)^2\right) + \exp\left(-N \left(\frac{\epsilon}{B_4}\right)^2\right) \right] - 4\delta_{\mathcal{C}}. \quad (184)$$

Now we lower bound the right side by replacing  $B_1, B_2, B_3, B_4$  with the upper bound  $\epsilon_1$  yielding

$$\mathbb{P}(\mathcal{E}_{good}(\epsilon)) \geq 1 - c_6 \exp\left(\log(\mathcal{C}_{\mathcal{F}_W}(\epsilon/\epsilon_2)) + \log\left(C_{\mathcal{F}_\theta}\left(\frac{\epsilon}{\epsilon_3}\right)\right) - c_7 N \left(\frac{\epsilon}{\epsilon_1}\right)^2\right) - 4\delta_{\mathcal{C}},$$

for some positive constants,  $c_6, c_7$ .

From Lemma 2 we have that for any  $\nu > 0$  it holds that

$$\log(\mathcal{C}_{\mathcal{F}_W}(\nu)) \leq R \log(\mathcal{C}_{\mathcal{F}_\theta}(L_\phi \nu / \gamma)). \quad (185)$$

Then we obtain

$$\mathbb{P}(\mathcal{E}_{good}(\epsilon)) \geq 1 - c_8 \exp\left(R \log\left(C_{\mathcal{F}_\theta}\left(\frac{L_\phi \epsilon}{\gamma \epsilon_2}\right)\right) - c_9 N \left(\frac{\epsilon}{\epsilon_1}\right)^2\right) - 4\delta_{\mathcal{C}}, \quad (186)$$

for some positive constants  $c_8, c_9$ .

From inequality (163), and (186) for any  $\epsilon \in [0, \epsilon_0]$  we have that

$$\begin{aligned} \mathbb{P}\left(|NC_\mu(\{W_j\}) - NC_{\mu_N}(\{W_j\})| \geq \lambda \Omega(f_\mu^*) \left[\Omega_{\mu_N}^\circ\left(-\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\}))\right) - 1\right] \right. \\ \left. - \frac{\alpha}{2} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 + (1 + \alpha)\epsilon + B_{eqI}(\mathcal{C}) + B_{plr}(\mathcal{C}) + (1 + \alpha)B_{nrm}(\mathcal{C})\right) \\ \leq c_8 \exp\left(R \log\left(C_{\mathcal{F}_\theta}\left(\frac{L_\phi \epsilon}{\gamma \epsilon_2}\right)\right) - c_9 N \left(\frac{\epsilon}{\epsilon_1}\right)^2\right) + 4\delta_{\mathcal{C}}. \end{aligned} \quad (187)$$

Next, we derive the operation conditions for  $\epsilon$  in terms of  $B_2, B_3$ , and  $B_4$ .

- $B_2 \geq B_3$ : observe that

$$B_2 \geq B_3 \iff \gamma \geq \Omega(f_\mu^*) L_\phi, \quad (188)$$

which establishes an upper bound on the regularization parameter.

- To establish a lower bound on regularization, we will require that  $\min\{B_2, B_3\} \geq B_4$ : We have that

$$\min\{B_2, B_3\} \geq B_4 \iff \quad (189)$$

$$\left(4 \min\left\{1, \frac{\Omega(f_\mu^*) L_\phi}{\gamma}\right\}\right) 4\gamma^2 \sigma_X^2 \sqrt{\left(1 + \|g\|_{\text{Lip}}^2 / \gamma^2\right) + \frac{\|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2}{\gamma^2 \sigma_X^2}} \geq 16\gamma^2 \sigma_X^2. \quad (190)$$

It is sufficient to have the below inequality to hold:

$$\min\left\{1, \frac{\Omega(f_\mu^*) L_\phi}{\gamma}\right\} \geq \frac{\gamma}{\sqrt{\gamma^2 + \|g\|_{\text{Lip}}^2}} \implies \min\{B_2, B_3\} \geq B_4. \quad (191)$$

Therefore,  $\gamma \geq \Omega(f_\mu^*) L_\phi$  is sufficient condition for  $B_2 \geq B_3 \geq B_4$ .

Then we have that

$$\epsilon_0 = \min\{B_1, B_4\} = 16n_Y \gamma^2 \sigma_X^2 \min\left\{1, \frac{L}{4} \left[1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left(1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2}\right)\right]\right\}, \quad (192)$$

and

$$\epsilon_1 = 16n_Y\gamma^2\sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}. \quad (193)$$

Now rescale the quantities  $\epsilon \leftarrow n_Y \frac{\epsilon}{(1+\alpha)}$ ,  $\epsilon_0 \leftarrow n_Y \epsilon_0$ , and  $\epsilon_1 \leftarrow n_Y \epsilon_1$ . Then we have

$$\epsilon_0 = 16\gamma^2\sigma_X^2 \min \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}; \quad (194)$$

$$\epsilon_1 = 16\gamma^2\sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}; \quad (195)$$

$$\epsilon_2 = \max\{8B_\ell\tilde{L}_\Phi, 8[\tilde{L}_\Phi B_\ell + \tilde{L}_\Phi B_\Phi L], 32\Omega(f_\mu^*)\tilde{L}_\Phi \max\{\tilde{L}_\Phi B_\ell/\tilde{L}_\Phi, LB_\Phi\}, 4\tilde{L}_\Phi B_\Phi\}, \quad (196)$$

and

$$\begin{aligned} \mathbb{P} \left( \frac{1}{n_Y} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \geq \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \right. \\ \left. - \frac{\alpha}{2n_Y} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 + \frac{1}{n_Y} [B_{eqI}(\mathcal{C}) + \lambda\Omega(f_\mu^*)B_{plr}(\mathcal{C}) + (1+\alpha)B_{nrm}(\mathcal{C})] + \epsilon \right) \\ \leq c_8 \exp \left( R \log \left( C_{\mathcal{F}_\theta} \left( \frac{L_\phi \epsilon}{\gamma \epsilon_2} \right) \right) - c_9 N \left( \frac{\epsilon}{(1+\alpha)\epsilon_1} \right)^2 \right) + 4\delta_{\mathcal{C}}. \end{aligned} \quad (197)$$

Now we bound the covering number under Assumption 4 and Lemma 2 via

$$\log \left( C_{\mathcal{F}_\theta} \left( \frac{L_\phi \epsilon}{\gamma \epsilon_2} \right) \right) \leq \dim(\mathcal{W}) \log(1 + 2\gamma\epsilon_2 r_\theta / (L_\phi \epsilon)) \leq c_{11} \dim(\mathcal{W}) \log(\gamma\epsilon_2 r_\theta / (L_\phi \epsilon)) \quad (198)$$

for some positive constant  $c_{11}$ .

Define

$$B(\mathcal{C}) := B_{eqI}(\mathcal{C}) + \lambda\Omega(f_\mu^*)B_{plr}(\mathcal{C}) + (1+\alpha)B_{nrm}(\mathcal{C}). \quad (199)$$

Then we have that

$$\begin{aligned} \mathbb{P} \left( \frac{1}{n_Y} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \geq \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \right. \\ \left. - \frac{\alpha}{2n_Y} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 + \frac{1}{n_Y} B(\mathcal{C}) + \epsilon \right) \\ \leq c_{12} \exp \left( c_{11} R \dim(\mathcal{W}) \log(\gamma\epsilon_2 r_\theta / (L_\phi \epsilon)) - c_{12} N \left( \frac{\epsilon}{(1+\alpha)\epsilon_1} \right)^2 \right) + 4\delta_{\mathcal{C}}. \end{aligned} \quad (200)$$

For some fixed  $\delta \in (0, 1]$  choose

$$\epsilon = \tilde{\Theta} \left( (1+\alpha)\epsilon_1 \sqrt{\frac{R \dim(\mathcal{W}) \log(\gamma\epsilon_2 r_\theta / L_\phi) \log(N) + \log(1/\delta)}{N}} \right). \quad (201)$$

Then the right side term of inequality (200) will be

$$\exp \left( R \dim(\mathcal{W}) \log(\gamma\epsilon_2 r_\theta / (L_\phi \epsilon)) - c_{12} N \left( \frac{\epsilon}{(1+\alpha)\epsilon_1} \right)^2 \right) = \tilde{\mathcal{O}}(\delta). \quad (202)$$

Rewriting the equation (200), we have

$$\begin{aligned}
 \mathbb{P} \left( \frac{1}{n_Y} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \gtrsim \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \right. \\
 \left. - \frac{\alpha}{2n_Y} \|f_\mu^* - \Phi_r(\{W_j\})\|_{\mu_N}^2 + \frac{1}{n_Y} B(\mathcal{C}) \right. \\
 \left. + (1 + \alpha) \epsilon_1 \sqrt{\frac{R \dim(\mathcal{W}) \log(\gamma \epsilon_2 r_\theta / L_\phi) \log(N) + \log(1/\delta)}{N}} \right) \lesssim \delta + \delta_{\mathcal{C}}. \quad (203)
 \end{aligned}$$

□

Theorem 4 has established generalization error for a generic parallel positively homogeneous network. Theorem 2 mentioned in the main text is a special case of Theorem 4, with the choice of convex set  $\mathcal{C} = \text{conv}(\mathcal{X} \times \mathbb{R}^{n_Y})$  by changing Assumption 7' to Assumption 7. Further,  $B(\mathcal{X} \times \mathbb{R}^{n_Y})$  will evaluate to 0, as  $\mathcal{P}_{\text{conv}(\mathcal{X} \times \mathbb{R}^{n_Y})}(\cdot)$  is just an identity operator.

## C APPLICATIONS

In this section, we apply our Theorem 4 for various applications. We apply our general theorem to low-rank matrix sensing, structured matrix sensing, two-layer linear neural network, two-layer ReLU neural network, and multi-head attention.

### C.1 Low-Rank Matrix Sensing

In this section, we state the corollary and its proof for matrix sensing, which is a direct consequence of Theorem 4. Firstly, we need to choose a convex set,  $\mathcal{C}$ , such that the Assumption 7' is satisfied. For matrix sensing we choose,  $\mathcal{C} = \{(X, \epsilon) : \|X\|_F \leq g, \|\epsilon\|_F \leq g\}$  to verify Assumption 7'. We need to compute,  $B(\mathcal{C})$ . This involves computing the expectation over the projection. Lemma 4 is pivotal for estimating  $B(\mathcal{C})$  in all the applications that are going to be discussed here.

**Lemma 4** (Projection of Gaussian vector on balls). *Consider a  $n$ -dimensional Gaussian vector  $\mathbf{x} \sim \mathcal{N}(0, (1/n)I_n)$ . Let  $M$  be a fixed matrix in  $\mathbb{R}^{n \times n}$  and  $\mathcal{A}$  be any set. Then*

$$\left| \langle M, \mathbb{E} [\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})\mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})^T] \mathbf{1}_{\mathcal{A}}(\mathbf{x}) \rangle \right| \leq \begin{cases} ge^{-g^2/2} \|M\|_2 & \text{if } g \geq 1 \\ \frac{1}{g} e^{-g^2/2} \|M\|_2 & \text{otherwise} \end{cases} \quad (204)$$

where  $\mathcal{P}_{\mathbb{B}(g)}(\cdot)$  is Euclidean projection onto the ball  $\mathbb{B}(g) := \{\mathbf{x} : \|\mathbf{x}\|_2 \leq g\}$ .

*Proof.* Define an event  $\mathcal{E} := \{\mathbf{x} \in \mathbb{B}(g)\}$ . When  $\mathcal{E}$  holds the function evaluates to zero,

$$\mathbb{E} [\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})\mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})^T] \mathbf{1}_{\mathcal{A}}(\mathbf{x}) = \mathbb{E} [(\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})\mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})^T) \mathbf{1}_{\mathcal{E}^c}(\mathbf{x}) \mathbf{1}_{\mathcal{A}}(\mathbf{x})], \quad (205)$$

so it suffices to consider the complement of the event  $\mathcal{E}$ . Now we take the inner product with  $M$  yielding

$$\begin{aligned}
 \left| \langle M, \mathbb{E} [\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})\mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})^T] \mathbf{1}_{\mathcal{A}}(\mathbf{x}) \rangle \right| &= \left| \langle M, \mathbb{E} [(\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})\mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})^T) \mathbf{1}_{\mathcal{E}^c}(\mathbf{x}) \mathbf{1}_{\mathcal{A}}(\mathbf{x})] \rangle \right| \\
 &\stackrel{(a)}{\leq} \|M\|_2 \mathbb{E} [(\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})\mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})^T) \mathbf{1}_{\mathcal{E}^c}(\mathbf{x}) \mathbf{1}_{\mathcal{A}}(\mathbf{x})] \| \\
 &\stackrel{(b)}{\leq} \|M\|_2 \mathbb{E} [\|(\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})\mathcal{P}_{\mathbb{B}(g)}(\mathbf{x})^T) \mathbf{1}_{\mathcal{E}^c}(\mathbf{x})\| \mathbf{1}_{\mathcal{A}}(\mathbf{x})] \\
 &\stackrel{(c)}{=} \|M\|_2 \mathbb{E} [\|\mathbf{x}\|_2^2 - g^2 | \mathbf{1}_{\mathcal{E}^c}(\mathbf{x}) \mathbf{1}_{\mathcal{A}}(\mathbf{x})] \\
 &\stackrel{(d)}{=} \|M\|_2 \mathbb{E} [|\|\mathbf{x}\|_2^2 - g^2| | \mathbf{1}_{\mathcal{E}^c}(\mathbf{x}) \mathbf{1}_{\mathcal{A}}(\mathbf{x})] \\
 &\stackrel{(e)}{\leq} \|M\|_2 \mathbb{E} [|\|\mathbf{x}\|_2^2 - g^2| | \mathbf{1}_{\mathcal{E}^c}(\mathbf{x})] \\
 &\stackrel{(f)}{=} \|M\|_2 \int_{\mathbf{x} \in \mathcal{E}^c} |\|\mathbf{x}\|_2^2 - g^2| \frac{1}{\sqrt{2\pi}} e^{-\frac{\|\mathbf{x}\|_2^2}{2}} d\mathbf{x} \\
 &\stackrel{(g)}{=} \|M\|_2 \left[ ge^{-g^2/2} - \sqrt{\frac{\pi}{2}} (g^2 - 1) \text{erfc}(g/\sqrt{2}) \right].
 \end{aligned}$$

The aforementioned computations involves (a) Cauchy-Schwartz inequality, (b) Jensen's inequality, (c) the norm of  $\mathcal{P}_{\mathcal{B}(g)}(\mathbf{x})$  when  $\mathbf{x} \in \mathcal{E}^c$  is  $g$ , (d) conditioning on indicator functions, (e) removing the conditioning increases the expectation over non-negative terms, (f) we apply the density of Gaussian, (g) standard normal integral. As a consequence of Theorem 1 from Zhang et al. (2020) we bound the complement error function,

$$\frac{e^{-z^2}}{\sqrt{\pi}z} \geq \text{erfc}(z) \geq \frac{2}{\sqrt{\pi}} \frac{e^{-z^2}}{z + \sqrt{z^2 + 2}}. \quad (206)$$

Then we have that

$$|\langle M, \mathbb{E} [[\mathbf{x}\mathbf{x}^T - \mathcal{P}_{\mathcal{B}(g)}(\mathbf{x})\mathcal{P}_{\mathcal{B}(g)}(\mathbf{x})^T] \mathbf{1}_{\mathcal{A}}(\mathbf{x})] \rangle| \leq \begin{cases} ge^{-g^2/2}\|M\|_2 & \text{if } g \geq 1 \\ \frac{1}{g}e^{-g^2/2}\|M\|_2 & \text{otherwise.} \end{cases} \quad (207)$$

□

Now, we state the generalization bound for the low-rank matrix sensing followed by its proof.

**Corollary 4** (Low-Rank Matrix Sensing). *Consider the true model for  $(X, y)$ , where  $X \in \mathbb{R}^{m \times n}$  is a random matrix with i.i.d. entries  $X_{lk} \sim \mathcal{N}(0, \frac{1}{mn})$  and  $y = \langle M^*, X \rangle + \epsilon$ , where  $M^* \in \mathbb{R}^{m \times n}$  and  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  is independent from  $X$ . For all  $i \in [N]$ , let  $(X_i, y_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{y} = \langle UV^T, X \rangle$ , where  $U \in \mathbb{R}^{m \times R}$  and  $V \in \mathbb{R}^{n \times R}$ . Let  $\delta \in (0, 1]$  be fixed. Define the non-convex problem*

$$\text{NC}_{\mu_N}^{\text{MS}}((U, V)) := \frac{1}{2N} \sum_{i=1}^N (y_i - \langle UV^T, X_i \rangle)^2 + \lambda \sum_{j=1}^R \|\mathbf{u}_j\|_2 \|\mathbf{v}_j\|_2, \quad (208)$$

and define  $\text{NC}_{\mu}^{\text{MS}}((U, V))$  similarly with the sum over  $i$  replaced by expectation taken over  $(X, y)$ .

Let  $(U, V)$  be a stationary point of  $\text{NC}_{\mu_N}^{\text{MS}}((U, V))$ . Suppose there exists  $C_{UV}, B_u, B_v > 0$  such that  $\|UV^T\|_2 \leq C_{UV}\|M^*\|_*$ , and for all  $j \in [R]$ ,  $\|\mathbf{u}_j\|_2 \leq B_u$ ,  $\|\mathbf{v}_j\|_2 \leq B_v$ . Then, with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned} \left| \text{NC}_{\mu}^{\text{MS}}((U, V)) - \text{NC}_{\mu_N}^{\text{MS}}((U, V)) \right| &\lesssim \|M^*\|_* \left[ \left\| \frac{1}{N} \sum_{i=1}^N (y_i - \langle UV^T, X_i \rangle) X_i \right\|_2 - \lambda \right] \\ &\quad + C_{UV}^2 \|M^*\|_*^2 \times \sqrt{\frac{R \log(R(C_{UV} + B_u B_v)) (m + n) \log(N) + \log(1/\delta)}{N}}. \end{aligned}$$

*Proof.* We set the following to obtain a generalization bound from Theorem 4 for the case of matrix sensing. First,

$$\ell(Y, \hat{Y}) = \frac{1}{2} \|Y - \hat{Y}\|_2^2 \implies (\alpha, L) = (0, 1); \quad (209)$$

$$\phi(W) = \langle \mathbf{u}\mathbf{v}^T, X \rangle; \quad (210)$$

$$\theta(W) = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2. \quad (211)$$

**Estimating  $\Omega(f_{\mu}^*)$ :** Since,  $M^*$  is the true matrix the regularizer at globally optimal solution can be upper bounded by Proposition 2,

$$\Omega(f_{\mu}^*) \leq \|M^*\|_*. \quad (212)$$

**Estimating  $\Omega_{\mu_N}^\circ(\cdot)$ :** Now we move on to compute the polar. We have

$$\begin{aligned}
 \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) &= \Omega_{\mu_N}^\circ \left( \frac{1}{\lambda} (g - \Phi_r(\{W_j\})) \right) \\
 &= \sup_{\|\mathbf{u}\| \leq 1; \|\mathbf{v}\| \leq 1} \frac{1}{N\lambda} \sum_{i=1}^N \langle Y_i - \langle UV^T, X_i \rangle, \mathbf{u}^T X_i \mathbf{v} \rangle \\
 &= \sup_{\|\mathbf{u}\| \leq 1; \|\mathbf{v}\| \leq 1} \frac{1}{N\lambda} \langle \mathbf{v}, \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle)^T \mathbf{u}^T X_i \rangle \\
 &= \sup_{\|\mathbf{u}\| \leq 1} \frac{1}{N\lambda} \left\| \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle) \mathbf{u}^T X_i \right\| \\
 &= \sup_{\|\mathbf{u}\| \leq 1} \frac{1}{N\lambda} \left\| \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle) X_i^T \mathbf{u} \right\| \\
 &= \frac{1}{\lambda} \left\| \frac{1}{N} \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle) X_i \right\|_2.
 \end{aligned}$$

**Defining  $\mathcal{F}_\theta$ :** Next, we estimate the relevant constants. First we estimate the constants from Assumption 4, suppose that  $\mathcal{B} := \{(\mathbf{u}, \mathbf{v}) : \|\mathbf{u}\|_2 \leq 1, \|\mathbf{v}\|_2 \leq 1\}$

$$\mathcal{F}_\theta := \{\mathbf{u}\mathbf{v}^T : \|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \leq 1\} \cap \mathcal{B} \quad (213)$$

**Estimating  $L_\phi$ :** The Lipschitz constant  $L_\phi$  in the function  $\mathcal{F}_\theta$  is  $L_\phi = \sup_{(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta} \|\langle \mathbf{u}\mathbf{v}^T, \cdot \rangle\|_{\text{Lip}} = 1$ .

**Estimating  $r_\theta$ :** We have that from A.M-G.M inequality,

$$\|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \leq \frac{1}{2} [\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2]. \quad (214)$$

Now for any  $(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta$  we have that  $0.5[\|\mathbf{u}\|_2 + \|\mathbf{v}\|_2] \leq 1$ . Therefore,  $\forall (\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta$  we have

$$\|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \leq \frac{1}{2} [\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2] \leq \sqrt{\frac{1}{2} [\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2]} \quad (215)$$

Then we need that  $\mathcal{F}_\theta \subseteq \mathbb{B}(r_\theta)$ , then must be  $r_\theta = \frac{1}{\sqrt{2}}$ .

**Defining  $\mathcal{F}_\mathcal{W}$ :** From the corollary's assumptions we have that,  $\mathcal{B}_R := \{(\mathbf{u}, \mathbf{v}) : \|\mathbf{u}\|_2 \leq B_u, \|\mathbf{v}\|_2 \leq B_v\}$ ; our hypothesis class is defined as

$$\mathcal{F}_\mathcal{W} := \{ \{(\mathbf{u}_j, \mathbf{v}_j)\} : \|\langle UV^T, \cdot \rangle\|_{\text{Lip}} = \|UV^T\|_2 \leq \gamma \} \cap \mathcal{B}_R. \quad (216)$$

As  $\gamma \geq \Omega(f_\mu^*) L_\phi$ , we choose  $\gamma = C_{UV} \|M^*\|_*$  for  $C_{UV} \geq 1$ . We have that

$$\mathcal{F}_\mathcal{W} = \{ \{(\mathbf{u}_j, \mathbf{v}_j)\} : \|\langle UV^T, \cdot \rangle\|_{\text{Lip}} = \|UV^T\|_2 \leq C_{UV} \|M^*\|_*, \|\mathbf{u}_j\| \leq B_u, \|\mathbf{v}_j\| \leq B_v \}. \quad (217)$$

**Estimating  $\epsilon_0$ :** From the data generating mechanism we have  $\|g\|_{\text{Lip}} = \|M^*\|_2$ ,  $\sigma_X = 1$ ,  $\sigma_{Y|X} = \sigma$ . Then we have the following constants from Theorem 4:

$$\epsilon_0 = 16\gamma^2 \sigma_X^2 \min \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (218)$$

which evaluates to

$$\epsilon_0 = 16C_{UV}^2 \|M^*\|_*^2 \min \left\{ 1, \frac{1 + \sigma^2}{4C_{UV}^2} \right\}. \quad (219)$$

From the corollary assumption we have that  $C_{UV} \leq 0.5\sqrt{1+\sigma^2}$ , which implies that

$$\epsilon_0 = 4(1 + \sigma^2)\|M^*\|_*^2. \quad (220)$$

**Estimating  $\epsilon_1$ :** Similarly, we evaluate

$$\epsilon_1 = 16\gamma^2\sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (221)$$

obtaining

$$\epsilon_1 = 16C_{UV}^2\|M^*\|_*^2 \max \left\{ 1, \frac{1 + \sigma^2}{4C_{UV}^2} \right\}. \quad (222)$$

From corollary assumption we have that  $C_{UV} \leq 0.5\sqrt{1+\sigma^2}$  which gives

$$\epsilon_1 = 16C_{UV}^2\|M^*\|_*^2. \quad (223)$$

**Defining convex set  $\mathcal{C}$ :** Consider a convex set  $\mathcal{C} = \mathbb{B}(g) = \{X : \|\text{vec}(X)\|_2 \leq g\}$ .

First and foremost we need to estimate  $\delta_{\mathcal{C}}$  for the following inequality to hold:

$$P(\cap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}. \quad (224)$$

The probability of  $X \in \mathcal{C} = \mathbb{B}(g)$  is equivalent to saying the probability of the event when  $\|\text{vec}(X)\|_2 \leq g$ . Since,  $X_{ij} \sim \mathcal{N}(0, 1/(m \times n))$  as a consequence of Bernstein's Inequality (Vershynin, 2018, Corollary 2.8.3) we have that for any  $t \geq 0$ ,

$$P(|\|\text{vec}(X)\|_2 - 1| \leq t) \geq 1 - 2 \exp(-cn_X t^2) \quad (225)$$

for some constant  $c \geq 0$ . Now we have

$$P(\|\text{vec}(X)\|_2 \leq g) \begin{cases} \geq 1 - 2 \exp(-cn_X(g-1)^2) & \text{if } g \geq 1 \\ \leq 2 \exp(-cn_X(g-1)^2) & \text{otherwise.} \end{cases} \quad (226)$$

We consider the case where  $g \geq 1$ , then we have that

$$P(\cap_{i=1}^N X_i \in \mathcal{C}) = P(\cap_{i=1}^N \|\text{vec}(X_i)\|_2 \leq g) \geq 1 - \underbrace{2N \exp(-cn_X(g-1)^2)}_{=\delta_{\mathcal{C}}}. \quad (227)$$

We have that  $\delta_{\mathcal{C}} = 2N \exp(-cn_X(g-1)^2)$ .

Now we evaluate  $B_{\ell}, B_{\Phi}, \tilde{L}_{\Phi}, \tilde{L}_{\phi}$ .

**Estimating  $B_{\Phi}$ :** Recall that  $r_{\theta} = \frac{1}{\sqrt{2}}$ . Then we have

$$B_{\Phi} = \sup_{Z \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|\langle UV^T, Z \rangle\| \quad (228)$$

$$= \sup_{Z \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|\langle \text{vec}(UV^T), \text{vec}(Z) \rangle\| \quad (229)$$

$$= g \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|\text{vec}(UV^T)\|_2 \quad (230)$$

$$= g \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \left\| \sum_{j=1}^R \text{vec}(\mathbf{u}_j \mathbf{v}_j^T) \right\|_2 \quad (231)$$

$$= gR \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|\text{vec}(\mathbf{u}_j \mathbf{v}_j^T)\|_2 \quad (232)$$

$$= gR \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|\mathbf{u}_j \mathbf{v}_j^T\|_F \quad (233)$$

$$= gR \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|\mathbf{u}_j\|_2 \|\mathbf{v}_j\|_2 \quad (234)$$

$$= gB_u B_v R. \quad (235)$$

**Estimating  $B_\ell$ :** Similarly, we have

$$B_\ell = \sup_{Z \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_\mathcal{W}} \|\langle UV^T - M^*, Z \rangle\| \quad (236)$$

$$= g \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_\mathcal{W}} \|\text{vec}(UV^T - M^*)\|_2 \quad (237)$$

$$\leq g[\|M^*\|_F + B_u B_v R]. \quad (238)$$

**Estimating  $\tilde{L}_\Phi$ :** Now, we compute the Lipschitz constant with respect to  $U, V$ . We have that

$$\tilde{L}_\Phi = \sup_{Z \in \mathcal{C}, (U, V), (U', V') \in \mathcal{F}_\mathcal{W}} \frac{\|\langle UV^T - U'V'^T, Z \rangle\|}{\max_j \sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (239)$$

$$= g \sup_{(U, V), (U', V') \in \mathcal{F}_\mathcal{W}} \frac{\|UV^T - U'V'^T\|_F}{\max_j \sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (240)$$

$$= gR \sup_{(U, V), (U', V') \in \mathcal{F}_\mathcal{W}} \frac{\|\mathbf{u}_j \mathbf{v}_j^T - \mathbf{u}'_j \mathbf{v}'_j{}^T\|_F}{\sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (241)$$

$$= gR \sup_{(U, V), (U', V') \in \mathcal{F}_\mathcal{W}} \frac{\|(\mathbf{u}_j - \mathbf{u}'_j) \mathbf{v}_j^T - \mathbf{u}'_j (\mathbf{v}_j - \mathbf{v}'_j)^T\|_F}{\sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (242)$$

$$\leq gR \sup_{(U, V), (U', V') \in \mathcal{F}_\mathcal{W}} \frac{\|(\mathbf{u}_j - \mathbf{u}'_j)\|_2 \|\mathbf{v}_j\|_2 + \|\mathbf{u}'_j\|_2 \|(\mathbf{v}_j - \mathbf{v}'_j)\|_2}{\sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (243)$$

$$\leq gR \sup_{(U, V), (U', V') \in \mathcal{F}_\mathcal{W}} \sqrt{\|\mathbf{v}_j\|_2^2 + \|\mathbf{u}'_j\|_2^2} \quad (244)$$

$$= g\sqrt{B_u^2 + B_v^2} R. \quad (245)$$

**Estimating  $\tilde{L}_\phi$ :** Similarly we get  $\tilde{L}_\phi = g\sqrt{B_u^2 + B_v^2}$ .

**Estimating  $\epsilon_2$ :** Recall that

$$\epsilon_2 = \max\{8B_\ell \tilde{L}_\Phi, 8\tilde{L}_\Phi[B_\ell + B_\Phi L], 32\Omega(f_\mu^*) \tilde{L}_\phi \max\{B_\ell, LB_\Phi\}, 4\tilde{L}_\Phi B_\Phi\}. \quad (246)$$

From all the constants computed earlier, we have that

$$\epsilon_2 = k_1 g^2 R^2 (\|M^*\|_F + B_u B_v)^2 \quad (247)$$

for some constant  $k_1 \geq 0$ .

Next we move on estimating  $B(\mathcal{C})$  we need to analyze three terms:

**The First Term:** We define the first term via

$$T_1 := \sup_{\{W_j\} \in \mathcal{F}_\mathcal{W}} \left| \|f_\mu^* \circ \mathcal{P}_\mathcal{C} - \Phi_r(\{W_j\}) \circ \mathcal{P}_\mathcal{C}\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right|.$$

For fixed  $(U, V)$ , we have

$$\begin{aligned} & \left| \|f_\mu^* \circ \mathcal{P}_\mathcal{C} - \Phi_r(\{W_j\}) \circ \mathcal{P}_\mathcal{C}\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \\ &= \left| \mathbb{E} [\langle M^* - UV^T, \mathcal{P}_\mathcal{C}(X) \rangle^2 - \langle M^* - UV^T, X \rangle^2] \right| \\ &= \left| \mathbb{E} \left[ \langle \text{vec}(M^* - UV^T) \text{vec}(M^* - UV^T)^T, \text{vec}(\mathcal{P}_\mathcal{C}(X)) \text{vec}(\mathcal{P}_\mathcal{C}(X))^T \rangle \right. \right. \\ &\quad \left. \left. - \langle \text{vec}(M^* - UV^T) \text{vec}(M^* - UV^T)^T, \text{vec}(X) \text{vec}(X)^T \rangle \right] \right| \\ &= \left| \langle \text{vec}(M^* - UV^T) \text{vec}(M^* - UV^T)^T, \mathbb{E} [\text{vec}(\mathcal{P}_\mathcal{C}(X)) \text{vec}(\mathcal{P}_\mathcal{C}(X))^T - \text{vec}(X) \text{vec}(X)^T] \rangle \right| \end{aligned}$$



From Lemma 4, taking  $g \geq 1$ ,

$$\begin{aligned} & \left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \\ & \leq ge^{-g^2/2} \|\text{vec}(M^* - UV^T) \text{vec}(M^* - UV^T)^T\|_2, \end{aligned} \quad (248)$$

whereupon further simplifying, we obtain

$$\left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \leq ge^{-g^2/2} \|M^* - UV^T\|_F^2. \quad (249)$$

Now, applying triangular inequality and taking the supremum, we obtain

$$T_1 \leq ge^{-g^2/2} (\|M^*\|_F + RB_u B_v)^2. \quad (250)$$

**The Second Term:** We define the second term via

$$\begin{aligned} T_2 := & \sup_{\{W_j\} \in \mathcal{F}_W, W' \in \mathcal{F}_\theta} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu \right. \\ & \left. - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right|. \end{aligned} \quad (251)$$

We have

$$\begin{aligned} & \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right| \\ & = |\mathbb{E} [\langle UV^T - M^*, \mathcal{P}_C(X) \rangle \langle \mathbf{u}\mathbf{v}^T, \mathcal{P}_C(X) \rangle - \langle UV^T - M^*, X \rangle \langle \mathbf{u}\mathbf{v}^T, X \rangle]| \\ & = |\langle \text{vec}(M^* - UV^T) \text{vec}(\mathbf{u}\mathbf{v}^T)^T, \mathbb{E} [\text{vec}(X) \text{vec}(X)^T - \text{vec}(\mathcal{P}_C(X)) \text{vec}(\mathcal{P}_C(X))^T] \rangle|. \end{aligned}$$

As a consequence of Lemma 4 we have

$$\begin{aligned} & \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right| \\ & \leq ge^{-g^2/2} \|\text{vec}(M^* - UV^T) \text{vec}(\mathbf{u}\mathbf{v}^T)^T\|_2 = ge^{-g^2/2} \|M^* - UV^T\|_F \|\mathbf{u}\mathbf{v}^T\|_F. \end{aligned} \quad (252)$$

Now we apply supremum over  $(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta$  and then  $(U, V)$  obtaining

$$T_2 \leq ge^{-g^2/2} [\|M^*\|_F + RB_u B_v]. \quad (253)$$

**The Third Term:** We define

$$\begin{aligned} T_3 := & \sup_{\{W_j\} \in \mathcal{F}_W} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \Phi_r(\{W_j\}) \circ \mathcal{P}_C \rangle_\mu \right. \\ & \left. - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu \right|. \end{aligned} \quad (254)$$

Similarly to the earlier item, we rewrite the above as

$$\left| \langle \text{vec}(M^* - UV^T) \text{vec}(UV^T)^T, \mathbb{E} [\text{vec}(X) \text{vec}(X)^T - \text{vec}(\mathcal{P}_C(X)) \text{vec}(\mathcal{P}_C(X))^T] \rangle \right| \quad (255)$$

As a consequence of Lemma 4 we have

$$\begin{aligned} & \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \Phi_r(\{W_j\}) \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu \right| \\ & \leq ge^{-g^2/2} \|\text{vec}(M^* - UV^T) \text{vec}(UV^T)^T\|_2 = ge^{-g^2/2} \|M^* - UV^T\|_F \|UV^T\|_F. \end{aligned} \quad (256)$$

Finally, we apply supremum over  $(U, V) \in \mathcal{F}_W$ , obtaining

$$T_3 \leq ge^{-g^2/2} B_u B_v R [\|M^*\|_F + RB_u B_v]. \quad (257)$$

Now combining equations (250), (253), (257) we obtain that

$$B(\mathcal{C}) \leq ge^{-g^2/2} [\alpha(\|M^*\|_F + RB_u B_v)^2 + \|M^*\|_F + RB_u B_v + B_u B_v R [\|M^*\|_F + RB_u B_v]] \quad (258)$$

We further upper bound for simplicity as

$$B(\mathcal{C}) \leq 4ge^{-g^2/2}(\|M^*\|_F + RB_u B_v)^2. \quad (259)$$

From Theorem 4 we have that

$$\frac{1}{n_Y} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \lesssim \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \quad (260)$$

$$+ \frac{4}{n_Y} ge^{-g^2/2} \{\|M^*\|_F + RB_u B_v\}^2 + 16C_{UV}^2 \|M^*\|_*^2 \times \left( \quad (261)$$

$$\sqrt{\frac{R \dim(\mathcal{W}) \log \left( C_{UV} \|M^*\|_* k_1 g^2 R^2 (\|M^*\|_F + B_u B_v)^2 \frac{1}{\sqrt{2}} \right) \log(N) + \log(1/\delta)}{N}} \right) \quad (262)$$

holds true w.p at least  $1 - \delta - 2N \exp(-cn_X(g-1)^2)$ .

Now choose

$$g = 1 + \mathcal{O} \left( \sqrt{\log(\sqrt{NR^{100}}) + \log(1/\delta)} \right). \quad (263)$$

Then we get that

$$\frac{1}{n_Y} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \quad (264)$$

$$\lesssim \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \quad (265)$$

$$+ \frac{4}{n_Y} \{\|M^*\|_F + RB_u B_v\}^2 \frac{\delta \sqrt{\log(NR^{100}) + \log(1/\delta)}}{NR^{100}} \quad (266)$$

$$+ 16C_{UV}^2 \|M^*\|_*^2 \sqrt{\frac{R(m+n) \log \left( C_{UV} \|M^*\|_* k_1 [\log(NR) + \log(1/\delta)] R^2 (\|M^*\|_F + B_u B_v)^2 \frac{1}{\sqrt{2}} \right) \log(N) + \log(1/\delta)}{N}} \quad (267)$$

holds true w.p at least  $1 - \delta$ . Now ignoring  $\log \log$  terms and keeping the right most term because of the dominance, we obtain

$$\begin{aligned} \frac{1}{n_Y} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| &\lesssim \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \\ &\quad + C_{UV}^2 \|M^*\|_*^2 \sqrt{\frac{R \log(R(C_{UV} + B_u B_v)) (m+n) \log(N) + \log(1/\delta)}{N}} \end{aligned} \quad (268)$$

holds true w.p at least  $1 - \delta$ .  $\square$

## C.2 Structured Matrix Sensing

Next, we move on to a slightly more generalized matrix sensing problem through which we impose certain structure in the factor  $U$ . Consider an atomic set  $\mathcal{U}$  that represents the set of structured columns, and suppose that  $U$  consists of columns that are affine combinations of the atoms in  $\mathcal{U}$ . We consider a gauge function  $\gamma_{\mathcal{U}}(\cdot)$  which is defined via

$$\gamma_{\mathcal{U}}(\mathbf{u}) := \inf \{t, t \geq 0 \text{ such that } \mathbf{u} \in t \text{conv}(\mathcal{U})\} \quad (269)$$

For instance,  $\mathcal{U}$  can be the intersection of  $L_2$  unit ball and  $L_1$  unit ball, which induces  $UV^T$  to be low-rank and  $U$  to be sparse. Imposing such structures has been well studied for convex problems by Chandrasekaran et al. (2012). Bach (2013) analyzed such structures for non-convex matrix factorization problems. However, their work was focused primarily on the optimization guarantees whereas our result below provides generalization/recovery guarantees for structured matrix sensing problems. We have the following corollary.

**Corollary 5** (Structured matrix sensing). *Consider the true model for  $(X, y)$ , where  $X \in \mathbb{R}^{m \times n}$  is a random matrix with i.i.d. entries  $X_{lk} \sim \mathcal{N}(0, \frac{1}{mn})$  and  $y = \langle U^* V^{*T}, X \rangle + \epsilon$ , where  $U^* \in \mathbb{R}^{m \times R^*}$ ,  $V^* \in \mathbb{R}^{n \times R^*}$  and  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  is independent from  $X$ . For all  $i \in [N]$ , let  $(X_i, y_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{y} = \langle UV^T, X \rangle$ , where  $U \in \mathbb{R}^{m \times R}$  and  $V \in \mathbb{R}^{n \times R}$ . Let  $\delta \in (0, 1]$  be fixed. Define the non-convex problem with the atomic set,  $\mathcal{U}$*

$$\text{NC}_{\mu_N}^{\text{SMS}}((U, V)) := \frac{1}{2N} \sum_{i=1}^N (y_i - \langle UV^T, X_i \rangle)^2 + \lambda \sum_{j=1}^R \gamma_{\mathcal{U}}(\mathbf{u}_j) \|\mathbf{v}_j\|_2, \quad (270)$$

and define  $\text{NC}_{\mu}^{\text{SMS}}((U, V))$  similarly with the sum over  $i$  replaced by expectation taken over  $(X, y)$ . Here  $\gamma_{\mathcal{U}}(\mathbf{u}) := \inf \{t; t \geq 0, \mathbf{u} \in t \text{conv}(\mathcal{U})\}$  for some specified atomic set,  $\mathcal{U}$ . Define

$$K_1 := \sum_{j=1}^{r^*} \gamma_{\mathcal{U}}(\mathbf{u}_j^*) \|\mathbf{v}_j^*\|_2; K_2 := \sup_{\|\mathbf{u}\| \leq 1} \gamma_{\mathcal{U}}(\mathbf{u}). \quad (271)$$

Let  $(U, V)$  be a stationary point of  $\text{NC}_{\mu_N}^{\text{SMS}}((U, V))$ . Suppose there exists  $C_{UV}, B_u, B_v > 0$  such that  $\|UV^T\|_2 \leq C_{UV} K_1$ , and for all  $j \in [R]$ ,  $\|\mathbf{u}_j\|_2 \leq B_u$ ,  $\|\mathbf{v}_j\|_2 \leq B_v$ . Then, with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned} \left| \text{NC}_{\mu}^{\text{SMS}}((U, V)) - \text{NC}_{\mu_N}^{\text{SMS}}((U, V)) \right| &\lesssim K_1 \left[ K_2 \left\| \frac{1}{N} \sum_{i=1}^N (y_i - \langle UV^T, X_i \rangle) X_i \right\|_2 - \lambda \right] \\ &+ C_{UV}^2 K_1^2 \sqrt{\frac{R \log(R(C_{UV} + B_u B_v)) (m+n) \log(N) + \log(1/\delta)}{N}}. \end{aligned} \quad (272)$$

**Remarks:** Similar to matrix sensing, the sample complexity required for consistency is only that  $N \gtrsim R(m+n)$  up to logarithmic terms, assuming a global minimum is found. The sample complexity is similar to low-rank matrix sensing (ignoring the scale and logarithmic dependency). To the best of our knowledge, this problem has not been studied from a statistical perspective, and our sample complexities match the corresponding convex slightly structured matrix sensing of Kakade et al. (2008). Unlike low-rank matrix sensing, the main technical challenge is to compute the polar/supremum term in the optimization error. In general, such a computation is NP-hard when the atomic set  $\mathcal{U}$  has non-negative atoms (Hendrickx and Olshevsky, 2010).

*Proof.* The proof is similar to that of Corollary 1, except for the computation of the polar. Therefore, we only compute the polar.

**Estimating  $\Omega(f_{\mu}^*)$ :** Since  $M^*$  is the true matrix the globally optimal solution would be  $M^*$ ; therefore, from Proposition 2 we have

$$\Omega(f_{\mu}^*) \leq \left( \sum_{j=1}^{r^*} \gamma_{\mathcal{U}}(\mathbf{u}_j^*) \|\mathbf{v}_j^*\|_2 \right). \quad (273)$$

**Estimating  $\Omega_{\mu_N}^\circ(\cdot)$ :** Now we move on to compute the polar.

$$\begin{aligned}
 \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) &= \Omega_{\mu_N}^\circ \left( \frac{1}{\lambda} (g - \Phi_r(\{W_j\})) \right) \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) \\
 &= \Omega_{\mu_N}^\circ \left( \frac{1}{\lambda} (g - \Phi_r(\{W_j\})) \right) \\
 &= \sup_{\gamma_{\mathcal{U}}(\mathbf{u}) \leq 1; \|\mathbf{v}\| \leq 1} \frac{1}{N\lambda} \sum_{i=1}^N \langle Y_i - \langle UV^T, X_i \rangle, \mathbf{u}^T X_i \mathbf{v} \rangle \\
 &= \sup_{\gamma_{\mathcal{U}}(\mathbf{u}) \leq 1; \|\mathbf{v}\| \leq 1} \frac{1}{N\lambda} \langle \mathbf{v}, \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle)^T \mathbf{u}^T X_i \rangle \\
 &= \sup_{\gamma_{\mathcal{U}}(\mathbf{u}) \leq 1} \frac{1}{N\lambda} \left\| \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle) \mathbf{u}^T X_i \right\| \\
 &= \sup_{\gamma_{\mathcal{U}}(\mathbf{u}) \leq 1} \frac{1}{N\lambda} \left\| \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle) X_i^T \mathbf{u} \right\|.
 \end{aligned}$$

This yields

$$\Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) \leq \left[ \sup_{\|\mathbf{u}\| \leq 1} \gamma_{\mathcal{U}}(\mathbf{u}) \right] \frac{1}{N\lambda} \left\| \sum_{i=1}^N (Y_i - \langle UV^T, X_i \rangle) X_i^T \right\|_2. \quad (274)$$

The rest of the proof is the same as that of low-rank matrix sensing (see section C.1).  $\square$

### C.3 Two-Layer Linear NN

Next, we consider the closely related problem of 2-Layer Linear Neural Networks, which is essentially a multi-dimensional matrix sensing problem; this is also referred to as non-convex linear regression. In practice, this approach has seemed to have better linear convergence (Arora et al., 2019) and generalization capabilities (Allen-Zhu et al., 2019) than vanilla linear regression. Corollary 6 provides generalization error upper bounds.

**Corollary 6** (2-Layer Linear Neural Network). *Consider the true model for  $(\mathbf{x}, \mathbf{y})$ , where  $\mathbf{x} \sim \mathcal{N}(0, (1/n)I_n) \in \mathbb{R}^n$ ,  $\mathbf{y} = U^* V^{*T} \mathbf{x} + \epsilon$ , where  $U^* \in \mathbb{R}^{m \times R^*}$ ,  $V^* \in \mathbb{R}^{n \times R^*}$ , and  $\epsilon \sim \mathcal{N}(0, (\sigma^2/m)I_m) \in \mathbb{R}^m$  independent from  $\mathbf{x}$ . For all  $i \in [N]$ , let  $(\mathbf{x}_i, \mathbf{y}_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{\mathbf{y}} = UV^T \mathbf{x}$ , where  $U \in \mathbb{R}^{m \times R}$ ,  $V \in \mathbb{R}^{n \times R}$ . Let  $\delta \in (0, 1]$  be fixed. Define the non-convex problem*

$$\text{NC}_{\mu_N}^{2\text{LNN}}((U, V)) := \frac{1}{2N} \sum_{i=1}^N \|\mathbf{y}_i - U[V^T \mathbf{x}_i]_+\|_2^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2), \quad (275)$$

and define  $\text{NC}_{\mu}^{2\text{LNN}}((U, V))$  similarly with the sum over  $i$  replaced by expectation taken over  $(\mathbf{x}, \mathbf{y})$ .

Let  $(U, V)$  be a stationary point of  $\text{NC}_{\mu_N}^{2\text{LNN}}((U, V))$ . Suppose there exists  $C_{UV}, B_u, B_v > 0$  such that  $\|UV^T\|_2 \leq C_{UV} [\|U^*\|_F^2 + \|V^*\|_F^2]$ , and for all  $j \in [R]$ ,  $\|\mathbf{u}_j\|_2 \leq B_u$ ,  $\|\mathbf{v}_j\|_2 \leq B_v$ . Then, with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned}
 \frac{1}{m} \left| \text{NC}_{\mu}^{2\text{LNN}}((U, V)) - \text{NC}_{\mu_N}^{2\text{LNN}}((U, V)) \right| &\lesssim \frac{1}{2m} [\|U^*\|_F^2 + \|V^*\|_F^2] \left[ \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2 \|\mathbf{x}_i\|_2 - \lambda \right] \\
 C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2]^2 &\sqrt{\frac{R \log(R(C_{UV} + B_u^2 + B_v^2))(m+n) \log(N) + \log(1/\delta)}{N}}.
 \end{aligned}$$

Similar to matrix sensing, we require that  $N \gtrsim R(m+n)$ , with  $\frac{R(m+n)}{N} \rightarrow 0$  for consistency at a global minimum. This matches classical results for (convex) linear regression.

*Proof.* To obtain a generalization bound from Theorem 4 for this setting, we set the following problem parameters:

$$\ell(Y, \hat{Y}) = \frac{1}{2} \|Y - \hat{Y}\| \implies (\alpha, L) = (0, 1); \quad (276)$$

$$\phi(W) = \langle \mathbf{v}, \mathbf{x} \rangle \mathbf{u}; \quad (277)$$

$$\theta(W) = \frac{1}{2} [\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2]. \quad (278)$$

**Estimating  $\Omega(f_\mu^*)$ :** From Proposition 2 we have that

$$\Omega(f_\mu^*) \leq \frac{\|U^*\|_F^2 + \|V^*\|_F^2}{2} \quad (279)$$

**Choosing  $\mathcal{F}_\theta$ :**

$$\mathcal{F}_\theta := \{(\mathbf{u}, \mathbf{v}) : \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 \leq 2, \|\mathbf{u}\|_2 \leq 1, \|\mathbf{v}\|_2 \leq 1\}. \quad (280)$$

**Estimating  $L_\phi$ :**

$$L_\phi = \sup_{(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta} \|\mathbf{u}\mathbf{v}^T(\cdot)\|_{\text{Lip}} = \sup_{(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta} \|\mathbf{u}\mathbf{v}^T\|_2 = 1. \quad (281)$$

**Estimating  $r_\theta$ :** For any  $(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta$ , we have that,

$$\frac{\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2}{2} \leq \sqrt{\frac{\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2}{2}} \implies \mathcal{F}_\theta \subseteq \mathbb{B}(1/\sqrt{2}). \quad (282)$$

Then we have  $r_\theta = 1/\sqrt{2}$ .

**Choosing  $\mathcal{F}_W$ :**

$$\mathcal{F}_W := \{(U, V) : \|UV^T\|_2 \leq \gamma, \|\mathbf{u}_j\| \leq B_u, \|\mathbf{v}_j\| \leq B_v\}. \quad (283)$$

As  $\gamma \geq \Omega(f_\mu^*)L_\phi = \frac{\|U^*\|_F^2 + \|V^*\|_F^2}{2}$ , we may take  $\gamma = C_{UV} \left[ \frac{\|U^*\|_F^2 + \|V^*\|_F^2}{2} \right]$  for some  $C_{UV}$ . We have that

$$\mathcal{F}_W = \left\{ \{(\mathbf{u}_j, \mathbf{v}_j)\} : \|\langle UV^T, \cdot \rangle\|_{\text{Lip}} = \|UV^T\|_2 \leq C_{UV} \left[ \frac{\|U^*\|_F^2 + \|V^*\|_F^2}{2} \right], \|\mathbf{u}_j\| \leq B_u, \|\mathbf{v}_j\| \leq B_v \right\}. \quad (284)$$

**Estimating  $\epsilon_0$ :** From the data generating mechanism we have  $\|g\|_{\text{Lip}} = \|M^*\|_2$ ,  $\sigma_X = 1$ ,  $\sigma_{Y|X} = \sigma$ , which yields the following constants from Theorem 4:

$$\epsilon_0 = 16\gamma^2\sigma_X^2 \min \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (285)$$

which evaluates to when  $C_{UV} \leq 0.5\sqrt{(1 + \sigma^2)}$

$$\epsilon_0 = 8C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2] \min \left\{ 1, \frac{1 + \sigma^2}{4C_{UV}^2} \right\} = 2 [\|U^*\|_F^2 + \|V^*\|_F^2]. \quad (286)$$

**Estimating  $\epsilon_1$ :** Similarly, we evaluate

$$\epsilon_1 = 16\gamma^2\sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (287)$$

obtaining

$$\epsilon_1 = 8C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2]. \quad (288)$$

**Choosing the convex set  $\mathcal{C}$ :** Consider a convex set  $\mathcal{C} = \mathbb{B}(g) = \{X : \|\text{vec}(X)\|_2 \leq g\}$ .

First and foremost we need to estimate  $\delta_{\mathcal{C}}$  for the following inequality to hold:

$$P(\cap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}. \quad (289)$$

The probability of  $\mathbf{x} \in \mathcal{C} = \mathbb{B}(g)$  is equivalent to saying the probability of the event when  $\|\mathbf{x}\|_2 \leq g$ . Since,  $x_i \sim \mathcal{N}(0, 1/n)$  as a consequence of Bernstein's Inequality (Vershynin, 2018, Corollary 2.8.3) we have that for any  $t \geq 0$ ,

$$P(|\|\mathbf{x}\|_2 - 1| \leq t) \geq 1 - 2\exp(-cn_X t^2) \quad (290)$$

for some constant  $c \geq 0$ . Now we have

$$P(\|\mathbf{x}\|_2 \leq g) \begin{cases} \geq 1 - 2 \exp(-cn_X(g-1)^2) & \text{if } g \geq 1 \\ \leq 2 \exp(-cn_X(g-1)^2) & \text{otherwise} \end{cases} \quad (291)$$

We consider the case where  $g \geq 1$ , then we have that

$$P(\cap_{i=1}^N X_i \in \mathcal{C}) = P(\cap_{i=1}^N \|\mathbf{x}\|_2 \leq g) \geq 1 - \underbrace{2N \exp(-cn_X(g-1)^2)}_{=\delta_{\mathcal{C}}}. \quad (292)$$

We have that  $\delta_{\mathcal{C}} = 2N \exp(-cn(g-1)^2)$ .

Now we evaluate  $B_{\ell}, B_{\Phi}, \tilde{L}_{\Phi}, \tilde{L}_{\phi}$ .

**Estimating  $B_{\Phi}$ :** We have

$$B_{\Phi} = \sup_{\mathbf{z} \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|UV^T \mathbf{z}\| \quad (293)$$

$$= g \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|UV^T\|_2 \quad (294)$$

$$= g\gamma \quad (295)$$

**Estimating  $B_{\ell}$ :** Similarly, we have

$$B_{\ell} = \sup_{\mathbf{z} \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|(UV^T - U^*V^{*T})\mathbf{z}\| \quad (296)$$

$$= g \sup_{\{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|UV^T - U^*V^{*T}\|_2 \quad (297)$$

$$= g(\|U^*V^{*T}\|_2 + \gamma). \quad (298)$$

**Estimating  $\tilde{L}_{\Phi}$ :** Now, we compute the Lipschitz constant with respect to  $U, V$ . We have

$$\tilde{L}_{\Phi} = \sup_{\mathbf{z} \in \mathcal{C}, (U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|(UV^T - U'V'^T)\mathbf{z}\|}{\max_j \sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (299)$$

$$= g \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|UV^T - U'V'^T\|_2}{\max_j \sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (300)$$

$$\leq g \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|UV^T - U'V'^T\|_F}{\max_j \sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (301)$$

$$= g \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|\sum_{j=1}^R \mathbf{u}_j \mathbf{v}_j^T - \mathbf{u}'_j \mathbf{v}'_j{}^T\|_F}{\max_j \sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (302)$$

$$= gR \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|\mathbf{u}_j \mathbf{v}_j^T - \mathbf{u}'_j \mathbf{v}'_j{}^T\|_F}{\sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (303)$$

$$= gR \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|(\mathbf{u}_j - \mathbf{u}'_j) \mathbf{v}_j^T - \mathbf{u}'_j (\mathbf{v}'_j - \mathbf{v}_j)^T\|_2}{\sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (304)$$

$$\leq gR \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|(\mathbf{u}_j - \mathbf{u}'_j)\|_2 \|\mathbf{v}_j\|_2 + \|\mathbf{u}'_j\|_2 \|(\mathbf{v}'_j - \mathbf{v}_j)\|_2}{\sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (305)$$

$$= gR \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \sqrt{\|\mathbf{v}_j\|_2^2 + \|\mathbf{u}'_j\|_2^2} \quad (306)$$

$$= g\sqrt{B_u^2 + B_v^2} R. \quad (307)$$

**Estimating  $\tilde{L}_\phi$ :** Similarly we get  $\tilde{L}_\phi = g\sqrt{B_u^2 + B_v^2}$ .

**Estimating  $\epsilon_2$ :** Recall that

$$\epsilon_2 = \max\{8B_\ell\tilde{L}_\phi, 8\tilde{L}_\phi[B_\ell + B_\Phi L], 32\Omega(f_\mu^*)\tilde{L}_\phi \max\{B_\ell, LB_\Phi\}, 4\tilde{L}_\phi B_\Phi\}. \quad (308)$$

From all the constants computed earlier, we have that

$$\epsilon_2 = k_1 g^2 R^2 C_{UV}^2 (\|U^*\|_F^2 + \|V^*\|_F^2) \sqrt{B_u^2 + B_v^2} \quad (309)$$

for some constant  $k_1 \geq 0$ .

Next, we move on to estimating  $B(\mathcal{C})$ . We need to analyze three terms:

**The First Term:** Define

$$T_1 := \sup_{\{W_j\} \in \mathcal{F}_W} \left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right|$$

We have

$$\left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| = \quad (310)$$

$$= \left| \mathbb{E} \left[ \|(U^* V^{*T} - UV^T) \mathcal{P}_C(\mathbf{x})\|^2 - \|(U^* V^{*T} - UV^T) \mathbf{x}\|^2 \right] \right| \quad (311)$$

$$= \left| \mathbb{E} \left[ \langle (U^* V^{*T} - UV^T)(U^* V^{*T} - UV^T)^T, (\mathcal{P}_C(\mathbf{x}))(\mathcal{P}_C(\mathbf{x}))^T \rangle - \langle (U^* V^{*T} - UV^T)(U^* V^{*T} - UV^T)^T, \mathbf{x} \mathbf{x}^T \rangle \right] \right| \quad (312)$$

$$= \left| \langle (U^* V^{*T} - UV^T)(U^* V^{*T} - UV^T)^T, \mathbb{E} [(\mathcal{P}_C(\mathbf{x}))(\mathcal{P}_C(\mathbf{x}))^T - \mathbf{x} \mathbf{x}^T] \rangle \right| \quad (313)$$

From Lemma 4 we obtain that (taking  $g \geq 1$ )

$$\begin{aligned} & \left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \\ & \leq g e^{-g^2/2} \|(U^* V^{*T} - UV^T)(U^* V^{*T} - UV^T)^T\|_2, \end{aligned} \quad (314)$$

on further simplifying, we get

$$\left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \leq g e^{-g^2/2} \|U^* V^{*T} - UV^T\|_2^2. \quad (315)$$

Now applying triangular inequality and taking the supremum, we obtain

$$T_1 \leq g e^{-g^2/2} (\|U^* V^{*T}\|_2 + \gamma)^2, \quad (316)$$

**The Second Term:** Define

$$\begin{aligned} T_2 := \sup_{\{W_j\} \in \mathcal{F}_W, W' \in \mathcal{F}_\theta} & \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu \right. \\ & \left. - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right|. \end{aligned} \quad (317)$$

We have

$$\begin{aligned} & \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right| = \\ & \left| \mathbb{E} \left[ \langle (UV^T - U^* V^{*T}) \mathcal{P}_C(\mathbf{x}), \mathbf{u} \mathbf{v}^T \mathcal{P}_C(\mathbf{x}) \rangle - \langle (UV^T - U^* V^{*T}) \mathbf{x}, \mathbf{u} \mathbf{v}^T \mathbf{x} \rangle \right] \right| \end{aligned} \quad (318)$$

is the same as

$$= \left| \langle (U^* V^{*T} - UV^T)(\mathbf{u} \mathbf{v}^T)^T, \mathbb{E} [\mathbf{x} \mathbf{x}^T - (\mathcal{P}_C(\mathbf{x}))(\mathcal{P}_C(\mathbf{x}))^T] \rangle \right| \quad (319)$$

As a consequence of Lemma 4 we have

$$\left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right|$$

$$\leq ge^{-g^2/2} \|(U^*V^{*T} - UV^T)(\mathbf{u}\mathbf{v}^T)^T\|_2 = ge^{-g^2/2} \|U^*V^{*T} - UV^T\|_F \|\mathbf{u}\mathbf{v}^T\|_F. \quad (320)$$

Now we apply supremum over  $(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta$ , and then over  $(U, V) \in \mathcal{F}_\mathcal{W}$ , yielding

$$T_2 \leq ge^{-g^2/2} \left[ \|U^*V^{*T}\|_2 + \gamma \right]. \quad (321)$$

**The Third Term:** Define

$$T_3 := \sup_{\{W_j\} \in \mathcal{F}_\mathcal{W}} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_\mathcal{C}, \Phi_r(\{W_j\}) \circ \mathcal{P}_\mathcal{C}), \Phi_r(\{W_j\}) \circ \mathcal{P}_\mathcal{C} \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu \right|. \quad (322)$$

Similarly to the earlier item, we rewrite the above as

$$= \left| \langle (U^*V^{*T} - UV^T)(UV^T)^T, \mathbb{E}[\mathbf{x}\mathbf{x}^T - (\mathcal{P}_\mathcal{C}(\mathbf{x}))(\mathcal{P}_\mathcal{C}(\mathbf{x}))^T] \rangle \right| \quad (323)$$

As a consequence of Lemma 4 we have

$$\begin{aligned} & \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_\mathcal{C}, \Phi_r(\{W_j\}) \circ \mathcal{P}_\mathcal{C}), \Phi_r(\{W_j\}) \circ \mathcal{P}_\mathcal{C} \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu \right| \\ & \leq ge^{-g^2/2} \|(U^*V^{*T} - UV^T)(UV^T)^T\|_2 \leq ge^{-g^2/2} \|U^*V^{*T} - UV^T\|_2 \|UV^T\|_2. \end{aligned} \quad (324)$$

Finally, we apply supremum over  $(U, V) \in \mathcal{F}_\mathcal{W}$ , obtaining

$$T_3 \leq ge^{-g^2/2} \gamma \left[ \|U^*V^{*T}\|_2 + \gamma \right]. \quad (325)$$

Now combining equations (316), (321), (325) we obtain that

$$B(\mathcal{C}) \leq ge^{-g^2/2} \left[ \alpha(\|U^*V^{*T}\|_2 + \gamma)^2 + \|U^*V^{*T}\|_2 + \gamma + \gamma \left[ \|U^*V^{*T}\|_2 + \gamma \right] \right]. \quad (326)$$

We further upper bound for simplicity as

$$B(\mathcal{C}) \leq ge^{-g^2/2} (1 + \gamma) (\|U^*V^{*T}\|_2 + \gamma). \quad (327)$$

From Theorem 4 we have that

$$\begin{aligned} & \frac{1}{n_Y} |\mathbf{NC}_\mu(\{W_j\}) - \mathbf{NC}_{\mu_N}(\{W_j\})| \lesssim \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \\ & + \frac{4}{n_Y} ge^{-g^2/2} (1 + \gamma) (\|U^*V^{*T}\|_2 + \gamma) + 8C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2]^2 \times \left( \right. \\ & \left. \sqrt{\frac{R(m+n) \log \left( k_1 g^2 R^2 C_{UV}^2 (\|U^*\|_F^2 + \|V^*\|_F^2) \sqrt{B_u^2 + B_v^2} C_{UV} [\|U^*\|_F^2 + \|V^*\|_F^2] \right) \log(N) + \log(1/\delta)}{N}} \right) \end{aligned} \quad (328)$$

holds true w.p at least  $1 - \delta - 2N \exp(-cn_X(g-1)^2)$ .

Now choose

$$g = 1 + \mathcal{O} \left( \sqrt{\log(N) + \log(1/\delta)} \right). \quad (329)$$

Now ignoring  $\log \log$  terms and keep the right most term because of the dominance,

$$\begin{aligned} & \frac{1}{n_Y} |\mathbf{NC}_\mu(\{W_j\}) - \mathbf{NC}_{\mu_N}(\{W_j\})| \lesssim \frac{\lambda}{n_Y} \Omega(f_\mu^*) \left[ \Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) - 1 \right] \\ & + C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2]^2 \sqrt{\frac{R \log(R(C_{UV} + B_u^2 + B_v^2)) (m+n) \log(N) + \log(1/\delta)}{N}} \end{aligned} \quad (330)$$

holds true w.p at least  $1 - \delta$ .

□



#### C.4 Two-Layer ReLU NN

Next, we present and prove the generalization bound for the two-layer ReLU neural network. This is one step ahead of all the linear models that were discussed earlier. Similarly to the Gaussian projections discussed in matrix sensing, we discuss ReLU projection results that will be used in the main proof.

**Lemma 5** (ReLU projection 1). *Consider  $U_1, U_2 \in \mathbb{R}^{m \times r}$ ,  $V_1, V_2 \in \mathbb{R}^{n \times r}$ . Denote, convex set  $\mathcal{C} = \mathbb{B}(g)$  that is  $g$ -radius hyper sphere, then we have that*

$$\begin{aligned} & \left| \mathbb{E} \left[ \|U_1[V_1^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+ - U_2[V_2^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+\|^2 - \|U_1[V_1^T \mathbf{x}]_+ - U_2[V_2^T \mathbf{x}]_+\|^2 \right] \right| \\ & \leq 2ge^{-g^2/2} [\|U_1\|_F^2 \|V_1\|_F^2 + \|U_2\|_F^2 \|V_2\|_F^2]. \end{aligned} \quad (331)$$

*Proof.* First, we re-write

$$\|U_1[V_1^T \mathbf{x}]_+ - U_2[V_2^T \mathbf{x}]_+\|^2 = \left\| \sum_{j=1}^r \mathbf{u}_{j1} [\mathbf{v}_{j1}^T \mathbf{x}]_+ - \mathbf{u}_{j2} [\mathbf{v}_{j2}^T \mathbf{x}]_+ \right\|^2 \quad (332)$$

$$= \left\| \sum_{j=1}^r \mathbf{u}_{j1} \mathbf{v}_{j1}^T \mathbf{x} \mathbf{1}_{\mathbf{v}_{j1}^T \mathbf{x} \geq 0} - \mathbf{u}_{j2} \mathbf{v}_{j2}^T \mathbf{x} \mathbf{1}_{\mathbf{v}_{j2}^T \mathbf{x} \geq 0} \right\|^2 \quad (333)$$

$$\begin{aligned} &= \sum_{j=1}^r \sum_{l=1}^r \left[ \langle (\mathbf{u}_{j1} \mathbf{v}_{j1}^T)^T (\mathbf{u}_{j1} \mathbf{v}_{j1}^T), \mathbf{x} \mathbf{x}^T \mathbf{1}_{\mathbf{v}_{j1}^T \mathbf{x} \geq 0} \rangle + \langle (\mathbf{u}_{j2} \mathbf{v}_{j2}^T)^T (\mathbf{u}_{j2} \mathbf{v}_{j2}^T), \mathbf{x} \mathbf{x}^T \mathbf{1}_{\mathbf{v}_{j2}^T \mathbf{x} \geq 0} \rangle \right. \\ & \quad \left. - 2 \langle (\mathbf{u}_{j2} \mathbf{v}_{j2}^T)^T (\mathbf{u}_{j1} \mathbf{v}_{j1}^T), \mathbf{x} \mathbf{x}^T \mathbf{1}_{\mathbf{v}_{j1}^T \mathbf{x} \geq 0} \mathbf{1}_{\mathbf{v}_{j2}^T \mathbf{x} \geq 0} \rangle \right]. \end{aligned} \quad (334)$$

Note that  $\mathbf{1}_{\mathbf{v}^T \mathbf{x} > 0} = \mathbf{1}_{\mathbf{v}^T \mathcal{P}_{\mathcal{C}}(\mathbf{x}) > 0}$ . Similarly, we have

$$\begin{aligned} \|U_1[V_1^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+ - U_2[V_2^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+\|^2 &= \sum_{j=1}^r \sum_{l=1}^r \left[ \langle (\mathbf{u}_{j1} \mathbf{v}_{j1}^T)^T (\mathbf{u}_{j1} \mathbf{v}_{j1}^T), \mathcal{P}_{\mathcal{C}}(\mathbf{x}) \mathcal{P}_{\mathcal{C}}(\mathbf{x})^T \mathbf{1}_{\mathbf{v}_{j1}^T \mathbf{x} \geq 0} \rangle \right. \\ & \quad + \langle (\mathbf{u}_{j2} \mathbf{v}_{j2}^T)^T (\mathbf{u}_{j2} \mathbf{v}_{j2}^T), \mathcal{P}_{\mathcal{C}}(\mathbf{x}) \mathcal{P}_{\mathcal{C}}(\mathbf{x})^T \mathbf{1}_{\mathbf{v}_{j2}^T \mathbf{x} \geq 0} \rangle \\ & \quad \left. - 2 \langle (\mathbf{u}_{j2} \mathbf{v}_{j2}^T)^T (\mathbf{u}_{j1} \mathbf{v}_{j1}^T), \mathcal{P}_{\mathcal{C}}(\mathbf{x}) \mathcal{P}_{\mathcal{C}}(\mathbf{x})^T \mathbf{1}_{\mathbf{v}_{j1}^T \mathbf{x} \geq 0} \mathbf{1}_{\mathbf{v}_{j2}^T \mathbf{x} \geq 0} \rangle \right]. \end{aligned} \quad (335)$$

Now computing the difference between equations (334), and (335) we get

$$\begin{aligned} & \left| \mathbb{E} \left[ \|U_1[V_1^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+ - U_2[V_2^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+\|^2 - \|U_1[V_1^T \mathbf{x}]_+ - U_2[V_2^T \mathbf{x}]_+\|^2 \right] \right| = \\ & \left| \sum_{j=1}^r \sum_{l=1}^r \mathbb{E} \left[ \langle (\mathbf{u}_{j1} \mathbf{v}_{j1}^T)^T (\mathbf{u}_{j1} \mathbf{v}_{j1}^T), (\mathcal{P}_{\mathcal{C}}(\mathbf{x}) \mathcal{P}_{\mathcal{C}}(\mathbf{x})^T - \mathbf{x} \mathbf{x}^T) \mathbf{1}_{\mathbf{v}_{j1}^T \mathbf{x} \geq 0} \rangle \right. \right. \\ & \quad + \langle (\mathbf{u}_{j2} \mathbf{v}_{j2}^T)^T (\mathbf{u}_{j2} \mathbf{v}_{j2}^T), (\mathcal{P}_{\mathcal{C}}(\mathbf{x}) \mathcal{P}_{\mathcal{C}}(\mathbf{x})^T - \mathbf{x} \mathbf{x}^T) \mathbf{1}_{\mathbf{v}_{j2}^T \mathbf{x} \geq 0} \rangle \\ & \quad \left. \left. - 2 \langle (\mathbf{u}_{j2} \mathbf{v}_{j2}^T)^T (\mathbf{u}_{j1} \mathbf{v}_{j1}^T), (\mathcal{P}_{\mathcal{C}}(\mathbf{x}) \mathcal{P}_{\mathcal{C}}(\mathbf{x})^T - \mathbf{x} \mathbf{x}^T) \mathbf{1}_{\mathbf{v}_{j1}^T \mathbf{x} \geq 0} \mathbf{1}_{\mathbf{v}_{j2}^T \mathbf{x} \geq 0} \rangle \right] \right|. \end{aligned} \quad (336)$$

After applying Lemma 4 and the triangular inequality we obtain, for  $g \geq 1$ ,

$$\begin{aligned} & \left| \mathbb{E} \left[ \|U_1[V_1^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+ - U_2[V_2^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+\|^2 - \|U_1[V_1^T \mathbf{x}]_+ - U_2[V_2^T \mathbf{x}]_+\|^2 \right] \right| \\ & \stackrel{(a)}{\leq} ge^{-g^2/2} \sum_{j=1}^r \left[ \sum_{l=1}^r \|\mathbf{u}_{j1}\|^2 \|\mathbf{v}_{j1}\|^2 \right. \\ & \quad \left. + 2 \|\mathbf{u}_{j1}\| \|\mathbf{u}_{j2}\| \|\mathbf{v}_{j1}\| \|\mathbf{v}_{j2}\| + \|\mathbf{u}_{j2}\|^2 \|\mathbf{v}_{j2}\|^2 \right] \\ & \stackrel{(b)}{\leq} 2ge^{-g^2/2} \sum_{j=1}^r \sum_{l=1}^r \left[ \|\mathbf{u}_{j1}\|^2 \|\mathbf{v}_{j1}\|^2 + \|\mathbf{u}_{j2}\|^2 \|\mathbf{v}_{j2}\|^2 \right] \\ & \stackrel{(c)}{\leq} 2ge^{-g^2/2} [\|U_1\|_F^2 \|V_1\|_F^2 + \|U_2\|_F^2 \|V_2\|_F^2]. \end{aligned} \quad (337)$$

In (a) we apply triangular inequality and apply Lemma 4, (b) we use the identity that  $a^2 + 2ab + b^2 = (a + b)^2$ , and (c) we use the identity that  $(a^2c^2 + b^2d^2) \leq (a^2 + b^2)(c^2 + d^2)$ . This completes the proof.  $\square$

**Lemma 6** (ReLU projection 2). *Consider  $U_1, U_2 \in \mathbb{R}^{m \times r}$ ,  $V_1, V_2 \in \mathbb{R}^{n \times r}$ . Denote the convex set  $\mathcal{C} = \mathbb{B}(g)$ ; that is, the  $g$ -radius hyper sphere. Then we have that*

$$\begin{aligned} & \left| \mathbb{E} \left[ \langle U_1[V_1^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+ - U_2[V_2^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+, U_1'[V_1'^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+ - U_2'[V_2'^T \mathcal{P}_{\mathcal{C}}(\mathbf{x})]_+ \rangle \right. \right. \\ & \quad \left. \left. - \langle U_1[V_1^T \mathbf{x}]_+ - U_2[V_2^T \mathbf{x}]_+, U_1'[V_1'^T \mathbf{x}]_+ - U_2'[V_2'^T \mathbf{x}]_+ \rangle \right] \right| \\ & \leq 2ge^{-g^2/2} [\|U_1\|_F \|U_1'\|_F \|V_1\|_F \|V_1'\|_F + \|U_2\|_F \|U_2'\|_F \|V_2\|_F \|V_2'\|_F]. \end{aligned} \quad (338)$$

*Proof.* The proof is similar to the proof of Lemma 5.  $\square$

**Corollary 7** (Two-Layer ReLU Neural Network). *Consider the true model for  $(\mathbf{x}, \mathbf{y})$ , where  $\mathbf{x} \sim \mathcal{N}(0, (1/n)I_n) \in \mathbb{R}^n$ ,  $\mathbf{y} = U^*[V^{*T} \mathbf{x}]_+ + \epsilon$ , where  $U^* \in \mathbb{R}^{m \times R^*}$ ,  $V^* \in \mathbb{R}^{n \times R^*}$ , and  $\epsilon \sim \mathcal{N}(0, (\sigma^2/m)I_m) \in \mathbb{R}^m$  independent from  $\mathbf{x}$ . For all  $i \in [N]$ , let  $(\mathbf{x}_i, \mathbf{y}_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{\mathbf{y}} = U[V^T \mathbf{x}]_+$ , where  $U \in \mathbb{R}^{m \times R}$ ,  $V \in \mathbb{R}^{n \times R}$ . Let  $\delta \in (0, 1]$  be fixed. Define the non-convex problem*

$$\text{NC}_{\mu_N}^{\text{ReLU}}((U, V)) := \frac{1}{2N} \sum_{i=1}^N \|\mathbf{y}_i - U[V^T \mathbf{x}_i]_+\|_2^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2), \quad (339)$$

and define  $\text{NC}_{\mu}^{\text{ReLU}}((U, V))$  similarly with the sum over  $i$  replaced by expectation taken over  $(\mathbf{x}, \mathbf{y})$ .

Let  $(U, V)$  be a stationary point of  $\text{NC}_{\mu_N}^{\text{ReLU}}((U, V))$ . Suppose there exists  $C_{UV}, B_u, B_v > 0$  such that  $\|UV^T\|_2 \leq C_{UV} [\|U^*\|_F^2 + \|V^*\|_F^2]$ , and for all  $j \in [R]$ ,  $\|\mathbf{u}_j\|_2 \leq B_u$ ,  $\|\mathbf{v}_j\|_2 \leq B_v$ . Then with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned} & \frac{1}{m} \left| \text{NC}_{\mu}^{\text{ReLU}}((U, V)) - \text{NC}_{\mu_N}^{\text{ReLU}}((U, V)) \right| \lesssim \frac{1}{2m} [\|U^*\|_F^2 + \|V^*\|_F^2] \left[ \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2 \|\mathbf{x}_i\|_2 - \lambda \right] \\ & + C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2] \left[ \frac{R(m+n) \log(R(m+n)(C_{UV} + B_u^2 + B_v^2)) \log(N) + \log(1/\delta)}{N} \right]^{1/2}. \end{aligned}$$

*Proof.* To obtain a generalization bound from Theorem 4 for this setting, we set the following problem parameters:

$$\ell(Y, \hat{Y}) = \frac{1}{2} \|Y - \hat{Y}\| \implies (\alpha, L) = (0, 1) \quad (340)$$

$$\phi(W) = [\langle \mathbf{v}, \mathbf{x} \rangle]_+ \mathbf{u}; \quad (341)$$

$$\theta(W) = \frac{1}{2} [\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2]. \quad (342)$$

**Estimating  $\Omega(f_{\mu}^*)$ :** From Proposition 2 we have that

$$\Omega(f_{\mu}^*) \leq \frac{\|U^*\|_F^2 + \|V^*\|_F^2}{2} \quad (343)$$

**Choosing  $\mathcal{F}_{\theta}$ :**

$$\mathcal{F}_{\theta} := \{(\mathbf{u}, \mathbf{v}) : \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 \leq 2, \|\mathbf{u}\|_2 \leq 1, \|\mathbf{v}\|_2 \leq 1\}. \quad (344)$$

**Estimating  $L_{\phi}$ :** The Lipschitz constant  $L_{\phi}$  in the function  $\mathcal{F}_{\theta}$  is  $L_{\phi} := \sup_{\|\mathbf{u}\| \leq 1, \|\mathbf{v}\| \leq 1} \|\mathbf{u}[\mathbf{v}^T \cdot]_+\|_{\text{Lip}} \leq \sup_{\|\mathbf{u}\| \leq 1, \|\mathbf{v}\| \leq 1} \|\mathbf{u}\| \|\mathbf{v}\| = 1$ .

**Estimating  $r_{\theta}$ :** For any  $(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_{\theta}$ , we have that,

$$\frac{\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2}{2} \leq \sqrt{\frac{\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2}{2}} \implies \mathcal{F}_{\theta} \subseteq \mathbb{B}(1/\sqrt{2}). \quad (345)$$

Then we have  $r_\theta = 1/\sqrt{2}$ .

**Choosing  $\mathcal{F}_W$ :** From the corollary's assumptions we have that  $\mathcal{B}_R := \{(\mathbf{u}, \mathbf{v}) : \|\mathbf{u}\|_2 \leq B_u, \|\mathbf{v}\|_2 \leq B_v\}$ ; our hypothesis class is defined as

$$\mathcal{F}_W := \{(U, V) : \|U[V^T \cdot]_+ \|_{\text{Lip}} \leq \|U\|_2 \|V\|_2 \leq \gamma, \|\mathbf{u}_j\| \leq B_u, \|\mathbf{v}_j\| \leq B_v\}. \quad (346)$$

From Proposition 2, we have that,  $\Omega(f_\mu^*) \leq \frac{1}{2} [\|U^*\|_F^2 + \|V^*\|_F^2]$ . As we require  $\gamma \geq \Omega(f_\mu^*) L_\phi = \frac{1}{2} [\|U^*\|_F^2 + \|V^*\|_F^2]$ , we set  $\gamma = C_{UV} \frac{[\|U^*\|_F^2 + \|V^*\|_F^2]}{2}$ .

$$\begin{aligned} \mathcal{F}_W = \left\{ \{(\mathbf{u}_j, \mathbf{v}_j)\} : \|U[V^T \cdot]_+ \|_{\text{Lip}} = \|UV^T\|_2 \leq \frac{C_{UV}}{2} [\|U^*\|_F^2 + \|V^*\|_F^2], \right. \\ \left. \|\mathbf{u}_j\| \leq B_u, \|\mathbf{v}_j\| \leq B_v \right\}. \end{aligned} \quad (347)$$

**Estimating  $\Omega_{\mu_N}^\circ(\cdot)$ :** We have

$$\Omega_{\mu_N}^\circ \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) = \Omega_{\mu_N}^\circ \left( \frac{1}{\lambda} (g - \Phi_r(\{W_j\})) \right) \quad (348)$$

$$= \sup_{\|\mathbf{u}\| \leq 1, \|\mathbf{v}\| \leq 1} \frac{1}{N\lambda} \sum_{i=1}^N \langle Y_i - U[V^T \mathbf{x}_i]_+, \mathbf{u}[\mathbf{v}^T \mathbf{x}_i]_+ \rangle \quad (349)$$

$$= \sup_{\|\mathbf{v}\| \leq 1} \frac{1}{N\lambda} \sum_{i=1}^N [\mathbf{v}^T \mathbf{x}_i]_+ \|Y_i - \hat{Y}_i\|_2 \quad (350)$$

$$\leq \frac{1}{N\lambda} \sum_{i=1}^N \sup_{\|\mathbf{v}\| \leq 1} [\mathbf{v}^T \mathbf{x}_i]_+ \|Y_i - \hat{Y}_i\|_2 \quad (351)$$

$$= \frac{1}{N\lambda} \sum_{i=1}^N \|\mathbf{x}_i\|_2 \|Y_i - \hat{Y}_i\|_2. \quad (352)$$

**Estimating  $\epsilon_0$ :** From the data generating mechanism we have  $\|g\|_{\text{Lip}} \leq \|U^*\|_2 \|V^*\|_2 \leq \frac{1}{2} [\|U^*\|_F + \|V^*\|_F]$ ,  $\sigma_X = 1$ ,  $\sigma_{Y|X} = \sigma$ . Then we have the following constants from Theorem 4:

$$\epsilon_0 = 16\gamma^2 \sigma_X^2 \min \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (353)$$

which evaluates to

$$\epsilon_0 = 8C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2] \min \left\{ 1, \frac{(1 + \sigma^2)}{4C_{UV}^2} \right\}. \quad (354)$$

Let  $C_{UV} \leq 0.5\sqrt{1 + \sigma^2}$  then we have

$$\epsilon_0 = 2(1 + \sigma^2) [\|U^*\|_F^2 + \|V^*\|_F^2]. \quad (355)$$

**Estimating  $\epsilon_1$ :** Similarly,

$$\epsilon_1 = 16\gamma^2 \sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (356)$$

obtaining

$$\epsilon_1 = 8C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2]. \quad (357)$$

**Defining a convex set  $\mathcal{C}$ :** Consider a convex set  $\mathcal{C} = \mathbb{B}(g) = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq g\}$ .

First and foremost we need to estimate,  $\delta_{\mathcal{C}}$  for the following inequality to hold:

$$P(\cap_{i=1}^N \mathbf{x}_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}. \quad (358)$$

The probability of  $\mathbf{x} \in \mathcal{C} = \mathbb{B}(g)$  is equivalent to saying the probability of the event when  $\|\mathbf{x}\|_2 \leq g$ . Since,  $x_i \sim \mathcal{N}(0, 1/n)$  as a consequence of Bernstein's Inequality (Vershynin, 2018, Corollary 2.8.3) we have that, for any  $t \geq 0$ ,

$$P(|\|\mathbf{x}\|_2 - 1| \leq t) \geq 1 - 2 \exp(-cn_X t^2), \quad (359)$$

for some constant  $c \geq 0$ . Now we have

$$P(\|\mathbf{x}\|_2 \leq g) \begin{cases} \geq 1 - 2 \exp(-cn_X(g-1)^2) & \text{if } g \geq 1 \\ \leq 2 \exp(-cn_X(g-1)^2) & \text{otherwise.} \end{cases} \quad (360)$$

We consider the case where  $g \geq 1$  yielding

$$P(\cap_{i=1}^N \mathbf{x}_i \in \mathcal{C}) = P(\cap_{i=1}^N \|\mathbf{x}_i\|_2 \leq g) \geq 1 - \underbrace{2N \exp(-cn_X(g-1)^2)}_{=\delta_{\mathcal{C}}}. \quad (361)$$

We have that  $\delta_{\mathcal{C}} = 2N \exp(-cn(g-1)^2)$ .

Now we evaluate  $B_{\ell}, B_{\Phi}, \tilde{L}_{\Phi}, \tilde{L}_{\phi}$ .

**Estimating  $B_{\Phi}$ :** We have

$$B_{\Phi} = \sup_{\mathbf{z} \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|U[V^T \mathbf{z}]_+\|_2 \quad (362)$$

$$\leq \sup_{\mathbf{z} \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|U\|_2 \| [V^T \mathbf{z}]_+ \|_2 \quad (363)$$

$$\leq \sup_{\mathbf{z} \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|U\|_2 \|V\|_2 \|\mathbf{z}\|_2 \quad (364)$$

$$= g\gamma. \quad (365)$$

**Estimating  $B_{\ell}$ :** Similarly, we have

$$B_{\ell} = \sup_{\mathbf{z} \in \mathcal{C}, \{(\mathbf{u}_j, \mathbf{v}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|U[V^T \mathbf{z}]_+ - U^*[V^{*T} \mathbf{z}]\|_2 \quad (366)$$

$$= 2g\gamma \quad (367)$$

**Estimating  $\tilde{L}_{\Phi}$ :** Now, we compute the Lipschitz constant with respect to  $U, V$ . We have

$$\tilde{L}_{\Phi} = \sup_{\mathbf{z} \in \mathcal{C}, (U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|U[V^T \mathbf{z}]_+ - U'[V'^T \mathbf{z}]_+\|}{\max_j \sqrt{\|\mathbf{u}_j - \mathbf{u}'_j\|^2 + \|\mathbf{v}_j - \mathbf{v}'_j\|^2}} \quad (368)$$

$$= R \sup_{\mathbf{z} \in \mathcal{C}, (U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|\mathbf{u}[\mathbf{v}^T \mathbf{z}]_+ - \mathbf{u}'[\mathbf{v}'^T \mathbf{z}]_+\|}{\sqrt{\|\mathbf{u} - \mathbf{u}'\|^2 + \|\mathbf{v} - \mathbf{v}'\|^2}} \quad (369)$$

$$= R \sup_{\mathbf{z} \in \mathcal{C}, (U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|(\mathbf{u} - \mathbf{u}')[\mathbf{v}^T \mathbf{z}]_+ - \mathbf{u}'[\mathbf{v}'^T \mathbf{z}]_+ - [\mathbf{v}^T \mathbf{z}]_+\|}{\sqrt{\|\mathbf{u} - \mathbf{u}'\|^2 + \|\mathbf{v} - \mathbf{v}'\|^2}} \quad (370)$$

$$\leq R \sup_{\mathbf{z} \in \mathcal{C}, (U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|(\mathbf{u} - \mathbf{u}')[\mathbf{v}^T \mathbf{z}]_+\| + \|\mathbf{u}'[\mathbf{v}'^T \mathbf{z}]_+ - [\mathbf{v}^T \mathbf{z}]_+\|}{\sqrt{\|\mathbf{u} - \mathbf{u}'\|^2 + \|\mathbf{v} - \mathbf{v}'\|^2}} \quad (371)$$

$$\leq R \sup_{\mathbf{z} \in \mathcal{C}, (U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{\|(\mathbf{u} - \mathbf{u}')[\mathbf{v}^T \mathbf{z}]_+\| + \|\mathbf{u}'[\mathbf{v}'^T \mathbf{z}]_+ - [\mathbf{v}^T \mathbf{z}]_+\|}{\sqrt{\|\mathbf{u} - \mathbf{u}'\|^2 + \|\mathbf{v} - \mathbf{v}'\|^2}} \quad (372)$$

$$\leq gR \sup_{(U, V), (U', V') \in \mathcal{F}_{\mathcal{W}}} \frac{B_v \|\mathbf{u} - \mathbf{u}'\| + B_u \|\mathbf{v} - \mathbf{v}'\|}{\sqrt{\|\mathbf{u} - \mathbf{u}'\|^2 + \|\mathbf{v} - \mathbf{v}'\|^2}} \quad (373)$$

$$= g\sqrt{B_u^2 + B_v^2} R. \quad (374)$$

**Estimating  $\tilde{L}_\phi$ :** Similarly we get  $\tilde{L}_\phi = g\sqrt{B_u^2 + B_v^2}$ .

**Estimating  $\epsilon_2$ :** Recall that

$$\epsilon_2 = \max\{8B_\ell\tilde{L}_\Phi, 8\tilde{L}_\Phi[B_\ell + B_\Phi L], 32\Omega(f_\mu^*)\tilde{L}_\phi \max\{B_\ell, LB_\Phi\}, 4\tilde{L}_\Phi B_\Phi\}. \quad (375)$$

From all the constants computed earlier, we have that

$$\epsilon_2 = k_1 g^2 R^2 C_{UV} \sqrt{B_u^2 + B_v^2} [\|U^*\|_F^2 + \|V^*\|_F^2], \quad (376)$$

for some constant  $k_1 \geq 0$ .

Next, we move on to estimating  $B(\mathcal{C})$ . We need to analyze three terms:

**The First Term:** Define

$$T_1 := \sup_{\{W_j\} \in \mathcal{F}_W} \left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right|. \quad (377)$$

For a fixed  $(U, V)$ , we have

$$\begin{aligned} & \left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \\ &= \left| \mathbb{E} \left[ \|U^*[V^{*T} \mathcal{P}_C(\mathbf{x})]_+ - U[V^T \mathcal{P}_C(\mathbf{x})]_+\|^2 - \|U^*[V^{*T} \mathbf{x}]_+ - U[V^T \mathbf{x}]_+\|^2 \right] \right|. \end{aligned} \quad (378)$$

From Lemma 5 taking  $g \geq 1$  we have

$$\begin{aligned} & \left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \\ & \leq 2ge^{-g^2/2} [\|U^*\|_F^2 \|V^*\|_F^2 + \|U\|_F^2 \|V\|_F^2], \end{aligned} \quad (379)$$

whereupon further simplifying, we obtain

$$\begin{aligned} & \sup_{(U,V) \in \mathcal{F}_W} \left| \|f_\mu^* \circ \mathcal{P}_C - \Phi_r(\{W_j\}) \circ \mathcal{P}_C\|_\mu^2 - \|f_\mu^* - \Phi_r(\{W_j\})\|_\mu^2 \right| \\ & \leq 2ge^{-g^2/2} (\|U^*\|_F^2 \|V^*\|_F^2 + R^2 \gamma^2). \end{aligned} \quad (380)$$

Now, applying triangular inequality and taking the supremum we obtain

$$T_2 \leq 2ge^{-g^2/2} (\|U^*\|_F^2 \|V^*\|_F^2 + R^2 \gamma^2). \quad (381)$$

**The Second Term:** Define

$$\begin{aligned} T_2 := & \sup_{\{W_j\} \in \mathcal{F}_W, W' \in \mathcal{F}_\theta} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu \right. \\ & \left. - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right|. \end{aligned} \quad (382)$$

For a fixed  $(U, V)$ ,  $(\mathbf{u}, \mathbf{v})$  we have

$$\begin{aligned} & \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right| = \\ & \left| \mathbb{E} \left[ \langle U[V^T \mathcal{P}_C(\mathbf{x})]_+ - U^*[V^{*T} \mathcal{P}_C(\mathbf{x})]_+, \mathbf{u}[\mathbf{v}^T \mathcal{P}_C(\mathbf{x})]_+ \rangle \right. \right. \\ & \quad \left. \left. - \langle U[V^T \mathcal{P}_C(\mathbf{x})]_+ - U^*[V^{*T} \mathbf{x}]_+, \mathbf{u}[\mathbf{v}^T \mathbf{x}]_+ \rangle \right] \right|. \end{aligned} \quad (383)$$

As a consequence of Lemma 6 we have

$$\left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right|$$

$$\leq 2ge^{-g^2/2} [\|U\|_F \|V\|_F \|\mathbf{u}\| \|\mathbf{v}\| + \|U^*\|_F \|V^*\|_F \|\mathbf{u}\| \|\mathbf{v}\|]. \quad (384)$$

Now we apply supremum over  $(\mathbf{u}, \mathbf{v}) \in \mathcal{F}_\theta$ , obtaining

$$T_2 \leq 2ge^{-g^2/2} [\|U\|_F \|V\|_F + \|U^*\|_F \|V^*\|_F]. \quad (385)$$

Finally, we apply supremum over  $(U, V) \in \mathcal{F}_W$ , obtaining

$$T_2 \leq 2ge^{-g^2/2} [\|U^*\|_F \|V^*\|_F + R\gamma]. \quad (386)$$

**The Third Term:** Define

$$T_3 := \sup_{\{W_j\} \in \mathcal{F}_W} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \Phi_r(\{W_j\}) \circ \mathcal{P}_C \rangle_\mu \right. \\ \left. - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu \right|. \quad (387)$$

For a fixed  $(U, V)$  we can rewrite the above to

$$\left| \mathbb{E} \left[ \langle U[V^T \mathcal{P}_C(\mathbf{x})]_+ - U^*[V^{*T} \mathcal{P}_C(\mathbf{x})]_+, U[V^T \mathcal{P}_C(\mathbf{x})]_+ \rangle \right. \right. \\ \left. \left. - \langle U[V^T \mathcal{P}_C(\mathbf{x})]_+ - U^*[V^{*T} \mathbf{x}]_+, U[V^T \mathbf{x}]_+ \rangle \right] \right|. \quad (388)$$

As a consequence of Lemma 4 we have

$$\left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \Phi_r(\{W_j\}) \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu \right| \\ \leq 2ge^{-g^2/2} \|U\|_F \|V\|_F [\|U\|_F \|V\|_F + \|U^*\|_F \|V^*\|_F]. \quad (389)$$

Finally, we apply supremum over  $(U, V) \in \mathcal{F}_W$ , obtaining

$$T_3 \leq 2ge^{-g^2/2} R\gamma [\|U^*\|_F \|V^*\|_F + R\gamma]. \quad (390)$$

Now combining  $T_1, T_2$  and  $T_3$  from equations (381), (386), (390) we have

$$B(\mathcal{C}) \leq 2ge^{-g^2/2} [\alpha(\|U^*\|_F^2 \|V^*\|_F + R^2 \gamma^2) + \|U\|_F \|V\|_F + R\gamma + R\gamma [\|U^*\|_F \|V^*\|_F + R\gamma]]. \quad (391)$$

We further upper bound for simplicity via

$$B(\mathcal{C}) \leq 4Rge^{-g^2/2} \gamma [\|U\|_F \|V\|_F + \gamma]. \quad (392)$$

From Theorem 4 we have that

$$\frac{1}{m} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \lesssim \frac{\lambda}{2m} [\|U^*\|_F^2 + \|V^*\|_F^2] \left[ \sup_{\|\mathbf{v}\| \leq 1} \frac{1}{N} \sum_{i=1}^N [\mathbf{v}^T \mathbf{x}_i]_+ \|Y_i - U[V^T \mathbf{x}_i]_+\| - \lambda \right] \\ + \frac{2}{m} Rge^{-g^2/2} C_{UV} [\|U^*\|_F^2 + \|V^*\|_F^2] [\|U\|_F \|V\|_F + C_{UV} [\|U^*\|_F^2 + \|V^*\|_F^2]] + C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2]^2 \\ \sqrt{\frac{R(m+n) \log \left( k_1 g^2 R^2 C_{UV}^2 \sqrt{B_u^2 + B_v^2} [\|U^*\|_F^2 + \|V^*\|_F^2]^2 \right) \log(N) + \log(1/\delta)}{N}}, \quad (393)$$

holds true w.p at least  $1 - \delta - 2N \exp(-cn_X(g-1)^2)$ .

Now choose

$$g = 1 + \mathcal{O} \left( \sqrt{\log(NR) + \log(1/\delta)} \right). \quad (394)$$

After ignoring all the log-log terms and using only dominant terms, we have that

$$\begin{aligned} & \frac{1}{m} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| \lesssim \\ & \frac{\lambda}{2m} [\|U^*\|_F^2 + \|V^*\|_F^2] \left[ \sup_{\|\mathbf{v}\| \leq 1} \frac{1}{N} \sum_{i=1}^N [\mathbf{v}^T \mathbf{x}_i]_+ \|Y_i - U[V^T \mathbf{x}_i]_+\| - \lambda \right] \\ & + C_{UV}^2 [\|U^*\|_F^2 + \|V^*\|_F^2]^2 \sqrt{\frac{R(m+n) \log(R(m+n)(C_{UV} + B_u^2 + B_v^2)) \log(N) + \log(1/\delta)}{N}}, \end{aligned} \quad (395)$$

holds true w.p at least  $1 - \delta$ .  $\square$

### C.5 Multi-head Attention

Next, we move to applying Theorem 4 to the single-layer multi-head attention problem. We require similar Gaussian projections arguments onto convex sets are needed to be established. For this application, we require Gaussian projections onto softmax, which is analyzed through Lemma 7. First, we define the soft max operation,  $\sigma_t(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ .

$$[\sigma_t(\mathbf{u})]_i := \exp(tu_i) / \left( \sum_{j=1}^n \exp(tu_j) \right), \quad (396)$$

where  $t$  is called the temperature.

A discrete version of soft-max is known as hard-max that is defined as

$$[\sigma(\mathbf{u})]_i := \mathbf{1}_{u_i = \max_i u_i}. \quad (397)$$

Note that when  $t \rightarrow \infty$ ,  $\sigma_t(\mathbf{u}) \rightarrow \sigma(\mathbf{u})$ .

**Lemma 7** (Gaussian Softmax Projection). *Let  $X \in \mathbb{R}^{m \times n}$  and  $X_{ij} \sim \mathcal{N}(0, 1/(mn))$  be independent random variables. Suppose  $\sigma_t(\cdot)$  is a softmax with temperature,  $t$ , and  $M$  is fixed matrix in  $\mathbb{R}^{m \times n}$ . Consider a convex set  $\mathcal{C} = \{X = (\mathbf{x}_1, \dots, \mathbf{x}_n) : \forall j \in [n]; \|\mathbf{x}_j\| \leq g\}$  for  $g \geq 1$ . Then*

$$\begin{aligned} & \sup_{\mathbf{z} \in \mathcal{F}_\mathbf{z}} |\langle M, \mathbb{E} [X \sigma_t(X^T \mathbf{z})(X \sigma_t(X^T \mathbf{z}))^T - \mathcal{P}_\mathcal{C}(X) \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z})(\mathcal{P}_\mathcal{C}(X) \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z}))^T] \rangle| \\ & \leq c_1 m^2 \|M\|_F g \exp(-c_2 m g^2). \end{aligned} \quad (398)$$

for some positive constant,  $c_1, c_2$ .

*Proof.* Denote

$$T := \sup_{\mathbf{z} \in \mathcal{F}_\mathbf{z}} |\langle M, \mathbb{E} [X \sigma_t(X^T \mathbf{z})(X \sigma_t(X^T \mathbf{z}))^T - \mathcal{P}_\mathcal{C}(X) \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z})(\mathcal{P}_\mathcal{C}(X) \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z}))^T] \rangle|. \quad (399)$$

Firstly, we upper bound the earlier term by Cauchy-Schwartz inequality:

$$T \leq \|M\|_F \sup_{\mathbf{z} \in \mathcal{F}_\mathbf{z}} \|\mathbb{E} [X \sigma_t(X^T \mathbf{z})(X \sigma_t(X^T \mathbf{z}))^T - \mathcal{P}_\mathcal{C}(X) \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z})(\mathcal{P}_\mathcal{C}(X) \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z}))^T]\|_F. \quad (400)$$

Denote  $a_t(\mathbf{z}, i, j) = \sigma_t(X^T \mathbf{z})_i \sigma_t(X^T \mathbf{z})_j$  and  $\tilde{a}_t(\mathbf{z}, i, j) = \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z})_i \sigma_t(\mathcal{P}_\mathcal{C}(X)^T \mathbf{z})_j$ ,

$$T \leq \|M\|_F \sup_{\mathbf{z} \in \mathcal{F}_\mathbf{z}} \left\| \sum_{i=1}^T \sum_{j=1}^T \mathbb{E} [\mathbf{x}_i \mathbf{x}_j^T a_t(\mathbf{z}, i, j) - \mathcal{P}_\mathcal{C}(\mathbf{x}_i) \mathcal{P}_\mathcal{C}(\mathbf{x}_j)^T \tilde{a}_t(\mathbf{z}, i, j)] \right\|_F. \quad (401)$$

Now we apply triangular inequality,

$$T \leq \|M\|_F \sup_{\mathbf{z} \in \mathcal{F}_\mathbf{z}} \sum_{i=1}^T \sum_{j=1}^T \|\mathbb{E} [\mathbf{x}_i \mathbf{x}_j^T a_t(\mathbf{z}, i, j) - \mathcal{P}_\mathcal{C}(\mathbf{x}_i) \mathcal{P}_\mathcal{C}(\mathbf{x}_j)^T \tilde{a}_t(\mathbf{z}, i, j)]\|_F. \quad (402)$$

Now we apply Cauchy-Schwartz,

$$T \leq \|M\|_F \sum_{i=1}^m \sum_{j=1}^m \sup_{\mathbf{z} \in \mathcal{F}_{\mathbf{z}}} \|\mathbb{E} [\mathbf{x}_i \mathbf{x}_j^T a_t(\mathbf{z}, i, j) - \mathcal{P}_{\mathcal{C}}(\mathbf{x}_i) \mathcal{P}_{\mathcal{C}}(\mathbf{x}_j)^T \tilde{a}_t(\mathbf{z}, i, j)]\|_F. \quad (403)$$

We can upper bound the earlier term via taking a supremum over indices  $i, j \in [T]$  then we have

$$T \leq \|M\|_F m^2 \sup_{i,j} \sup_{\mathbf{z} \in \mathcal{F}_{\mathbf{z}}} \|\mathbb{E} [\mathbf{x}_i \mathbf{x}_j^T a_t(\mathbf{z}, i, j) - \mathcal{P}_{\mathcal{C}}(\mathbf{x}_i) \mathcal{P}_{\mathcal{C}}(\mathbf{x}_j)^T \tilde{a}_t(\mathbf{z}, i, j)]\|_F. \quad (404)$$

Observe that argument inside the expectation is 0 on the event  $X \in \mathcal{C}$ , thereby we only have the case where  $X \notin \mathcal{C}$  then we have

$$T \leq \|M\|_F m^2 \sup_{i,j} \sup_{\mathbf{z} \in \mathcal{F}_{\mathbf{z}}} \|\mathbb{E} [\mathbf{x}_i \mathbf{x}_j^T a_t(\mathbf{z}, i, j) - \mathcal{P}_{\mathcal{C}}(\mathbf{x}_i) \mathcal{P}_{\mathcal{C}}(\mathbf{x}_j)^T \tilde{a}_t(\mathbf{z}, i, j) \mathbf{1}_{\mathcal{E}^c}]\|_F. \quad (405)$$

On application Cauchy-Schwartz identity again we have

$$T \leq \|M\|_F m^2 \sup_{i,j} \sup_{\mathbf{z} \in \mathcal{F}_{\mathbf{z}}} \mathbb{E} [\|\mathbf{x}_i \mathbf{x}_j^T a_t(\mathbf{z}, i, j) - \mathcal{P}_{\mathcal{C}}(\mathbf{x}_i) \mathcal{P}_{\mathcal{C}}(\mathbf{x}_j)^T \tilde{a}_t(\mathbf{z}, i, j)\|_F \mathbf{1}_{\mathcal{E}^c}]. \quad (406)$$

We have that  $\mathbf{x} \notin \mathcal{C}, \mathcal{P}_{\mathcal{C}}(\mathbf{x}) = g\mathbf{x}/\|\mathbf{x}\|$ . By using this fact we have

$$T \leq \|M\|_F m^2 \sup_{i,j} \sup_{\mathbf{z} \in \mathcal{F}_{\mathbf{z}}} \mathbb{E} \left[ \left\| \left( a_t(\mathbf{z}, i, j) - \frac{g^2}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \tilde{a}_t(\mathbf{z}, i, j) \right) \mathbf{x}_i \mathbf{x}_j^T \right\|_F \mathbf{1}_{\mathcal{E}^c} \right]. \quad (407)$$

Now we recall Reverse Fatou's Lemma, for any function sequence,  $f_n \in L^2(\mu)$ , we have

$$\limsup_{n \rightarrow \infty} \int f_n d\mu \leq \int \limsup_{n \rightarrow \infty} f_n d\mu. \quad (408)$$

on applying this identity we have

$$T \leq \|M\|_F m^2 \sup_{i,j} \mathbb{E} \left[ \sup_{\mathbf{z} \in \mathcal{F}_{\mathbf{z}}} \left\| a_t(\mathbf{z}, i, j) - \frac{g^2}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \tilde{a}_t(\mathbf{z}, i, j) \right\| \|\mathbf{x}_i \mathbf{x}_j^T\|_F \mathbf{1}_{\mathcal{E}^c} \right]. \quad (409)$$

Observe that  $a_t(\mathbf{z}, i, j) \leq 1$  and  $\frac{g^2}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \tilde{a}_t(\mathbf{z}, i, j) \leq 1$  when  $X \notin \mathcal{C}$  therefore we have

$$T \leq \|M\|_F m^2 \sup_{i,j} \mathbb{E} [\|\mathbf{x}_i \mathbf{x}_j^T\|_F \mathbf{1}_{\mathcal{E}^c}]. \quad (410)$$

Now since  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are iid we have

$$T \leq \|M\|_F m^2 \sup_{i,j} \mathbb{E} [\|\mathbf{x}_i\|_2 \mathbf{1}_{\mathcal{E}^c}] \mathbb{E} [\|\mathbf{x}_j\|_2 \mathbf{1}_{\mathcal{E}^c}]. \quad (411)$$

By Gaussian integral over norm for  $g \geq 1$  we have

$$T \lesssim m^2 \|M\|_F g \exp(-mg^2) \quad (412)$$

□

With the above result on Gaussian softmax projection, we now state the corollary for the single-layer multi-head attention problem and its proof.

**Corollary 8** (Transformers). *Consider the true model for  $(X, \mathbf{y})$ , where  $X \in \mathbb{R}^{n \times T}$  is a random matrix with i.i.d. entries  $X_{lk} \sim \mathcal{N}(0, 1/(nT))$  and  $\mathbf{y} = A^* X \mathbf{b}^* + \epsilon$ , where  $A^* \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b}^* \in \mathbb{S}^{T-1}$  and  $\epsilon \sim \mathcal{N}(0, (\sigma^2/m)I_m)$  is independent from  $X$ . For all  $i \in [N]$ , let  $(X_i, \mathbf{y}_i)$  be i.i.d. samples from this true model. Consider the estimator  $\hat{\mathbf{y}} = \sum_{j=1}^R V_j X \sigma(X^T \mathbf{z}_j)$ ,  $V_j \in \mathbb{R}^n$ ,  $\mathbf{z}_j \in \mathbb{R}^n$ . Let  $\delta \in (0, 1]$  be fixed. Define the non-convex problem*



$$\text{NC}_{\mu_N}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\}) := \frac{1}{2N} \sum_{i=1}^N \|\mathbf{y}_i - \sum_{j=1}^R V_j X_i \sigma_t(X_i^T \mathbf{z}_j)\|_2^2 + \lambda \sum_{j=1}^R [\|V_j\|_F + \delta_{\{\mathbf{z}: \|\mathbf{z}\|_2 \leq 1\}}(\mathbf{z}_j)],$$

where,  $\sigma_t(\cdot)$  is softmax function with temperature  $t$ , for  $k \in [T]$  defined  $\sigma_t(\mathbf{u})_k := \exp(tu_k) / \sum_{l=1}^T \exp(tu_l)$  and define  $\text{NC}_{\mu}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\})$  similarly with the sum over  $i$  replaced by expectation taken over  $(X, \mathbf{y})$ .

Let  $\{(V_j, \mathbf{z}_j)\}$  be a stationary point of  $\text{NC}_{\mu_N}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\})$ . Suppose there exists  $C_V, B_V > 0$  such that  $\sum_{j=1}^R \|V_j\|_F \leq C_V \|A^*\|_F$ , and for all  $j \in [R]$ ,  $\|V_j\|_F \leq B_V$ .

Then with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned} \frac{1}{m} \left| \text{NC}_{\mu}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\}) - \text{NC}_{\mu_N}^{\text{TF}}(\{(V_j, \mathbf{z}_j)\}) \right| &\lesssim \frac{1}{2m} \|A^*\|_F \left[ \frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2 \|X_i\|_2 - \lambda \right] \\ &+ C_V^2 \|A^*\|_F^2 \sqrt{\frac{R(m+n) \log(R(m+n)(C_V+B_V)) \log(N) + \log(1/\delta)}{N}}. \end{aligned}$$

*Proof.* To obtain a generalization bound from Theorem 4 for the case of matrix sensing, we set the following problem parameters.

$$\ell(Y, \hat{Y}) = \frac{1}{2} \|Y - \hat{Y}\| \implies (\alpha, L) = (0, 1); \quad (413)$$

$$\phi(W) = V X \sigma_t(X^T \mathbf{z}); \quad (414)$$

$$\theta(W) = \|V\|_F + \delta_{\mathbf{z} \in \mathbb{B}(1)}. \quad (415)$$

**Estimating  $\Omega_{\mu_N}(\cdot)$ :** Now we move on to compute the polar:

$$\begin{aligned} \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) &= \Omega_{\mu_N}^{\circ} \left( \frac{1}{\lambda} (g - \Phi_r(\{W_j\})) \right), \\ &= \sup_{\|V\|_F \leq 1, \|\mathbf{z}\| \leq 1} \frac{1}{N\lambda} \sum_{i=1}^N \langle Y_i - \hat{Y}_i, V X_i \sigma_t(X_i^T \mathbf{z}) \rangle, \\ &= \sup_{\|V\|_F \leq 1, \|\mathbf{z}\| \leq 1} \frac{1}{N\lambda} \left\| \sum_{i=1}^N \langle (Y_i - \hat{Y}_i)(X_i \sigma_t(X_i^T \mathbf{z}))^T, V \rangle \right\|, \\ &= \sup_{\|V\|_F \leq 1, \|\mathbf{z}\| \leq 1} \frac{1}{N\lambda} \left\| \sum_{i=1}^N \langle (Y_i - \hat{Y}_i)(X_i \sigma_t(X_i^T \mathbf{z}))^T, V \rangle \right\|, \\ \Omega_{\mu_N}^{\circ} \left( -\frac{1}{\lambda} \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})) \right) &\leq \sup_{\|\mathbf{z}\| \leq 1} \frac{1}{N\lambda} \left\| \sum_{i=1}^N (Y_i - \hat{Y}_i)(X_i \sigma_t(X_i^T \mathbf{z}))^T \right\|_F, \\ &\leq \frac{1}{N\lambda} \sum_{i=1}^N \sup_{\|\mathbf{z}\| \leq 1} \|(Y_i - \hat{Y}_i)(X_i \sigma_t(X_i^T \mathbf{z}))^T\|_F, \\ &\leq \frac{1}{N\lambda} \sum_{i=1}^N \sup_{\|\mathbf{z}\| \leq 1} \|Y_i - \hat{Y}_i\|_F \|(X_i \sigma_t(X_i^T \mathbf{z}))^T\|_F, \\ &\leq \frac{1}{N\lambda} \sum_{i=1}^N \sup_{\|\mathbf{z}\| \leq 1} \|Y_i - \hat{Y}_i\|_F \|(X_i \sigma_t(X_i^T \mathbf{z}))^T\|_F, \\ &\leq \frac{1}{N\lambda} \sum_{i=1}^N \|Y_i - \hat{Y}_i\|_2 \|X_i\|_2. \end{aligned}$$

**Choose  $\mathcal{F}_\theta$ :** From Assumptions 4 suppose that

$$\mathcal{F}_\theta := \{(V, \mathbf{z}) : \|V\|_2 \leq 1, \|\mathbf{z}\|_2 \leq 1\}; \quad (416)$$

**Computing  $L_\phi$ :** The Lipschitz constant  $L_\phi$  in the function  $\mathcal{F}_\theta$  is  $L_\phi := \sup_{\|V\| \leq 1, \|\mathbf{z}\| \leq 1} \|V(\cdot)\sigma((\cdot)^T \mathbf{z})\|_{\text{Lip}} \leq \sup_{\|V\| \leq 1} \|V\| = 1$ .

**Computing  $r_\phi$ :** Clearly, when  $r_\theta = \sqrt{2}$  we have that  $\mathcal{F}_\theta \subseteq \mathbb{B}(\sqrt{2})$ .

**Choose  $\mathcal{F}_\mathcal{W}$ :** From the corollary's assumptions we have that,  $\mathcal{B}_R := \{(V, \mathbf{z}) : \|V\|_2 \leq B_V, \|\mathbf{z}\|_2 \leq 1\}$ ; our hypothesis class is defined as

$$\mathcal{F}_\mathcal{W} := \{(V_j, \mathbf{z}_j)\} : \left\| \sum_{j=1}^r V_j \cdot \sigma((\cdot)^T \mathbf{z}_j) \right\|_{\text{Lip}} \leq \sum_{j=1}^r \|V_j\|_F \leq \gamma, \|V_j\|_F \leq B_V, \|\mathbf{z}_j\| \leq 1\}. \quad (417)$$

From Proposition 2 we have  $\Omega(f_\mu^*) \leq \|A^*\|_F$ . We have  $\gamma \geq \Omega(f_\mu^*)L_\phi = \|A^*\|_F$ , then we set  $\gamma = C_V \|A^*\|_F$ . We have that

$$\begin{aligned} \mathcal{F}_\mathcal{W} := \left\{ \{(V_j, \mathbf{z}_j)\} : \left\| \sum_{j=1}^r V_j \cdot \sigma((\cdot)^T \mathbf{z}_j) \right\|_{\text{Lip}} \leq \sum_{j=1}^r \|V_j\|_F \leq C_V \|A^*\|_F, \right. \\ \left. \|V_j\|_F \leq B_V, \|\mathbf{z}_j\| \leq 1 \right\}. \end{aligned} \quad (418)$$

**Estimating  $\epsilon_0$ :** From the data generating mechanism we have  $\|g\|_{\text{Lip}} \leq \|A^*\|_F$ ,  $\sigma_X = 1$ ,  $\sigma_{Y|X} = \sigma$ , then we have the following constants from Theorem 4:

$$\epsilon_0 = 16\gamma^2 \sigma_X^2 \min \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (419)$$

this evaluates to,

$$\epsilon_0 = 16(\sigma^2 + 1) \|A^*\|_F^2, \quad (420)$$

when  $C_V \leq \sqrt{1 + \sigma^2}$ .

**Estimating  $\epsilon_1$ :** Similarly, we evaluate

$$\epsilon_1 = 16\gamma^2 \sigma_X^2 \max \left\{ 1, \frac{L}{4} \left[ 1 + \frac{\|g\|_{\text{Lip}}^2}{\gamma^2} \left( 1 + \frac{\sigma_{Y|X}^2}{\sigma_X^2} \right) \right] \right\}, \quad (421)$$

obtaining,

$$\epsilon_1 = 16C_V^2 \|A^*\|_F^2, \quad (422)$$

when  $C_V \leq \sqrt{1 + \sigma^2}$ .

**Choosing the convex set  $\mathcal{C}$ :** Consider a convex set  $\mathcal{C} = \{X = (\mathbf{x}_1, \dots, \mathbf{x}_T) : \|\mathbf{x}_j\|_2 \leq g/\sqrt{T}\}$ .

First and foremost we need to estimate  $\delta_{\mathcal{C}}$  for the inequality to hold:

$$P(\cap_{i=1}^N \mathbf{x}_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}. \quad (423)$$

The probability of  $\mathbf{x} \in \mathcal{C} = \mathbb{B}(g)$  is equivalent to saying the probability of the event when  $\|\mathbf{x}\|_2 \leq g$ . Since,  $x_i \sim \mathcal{N}(0, 1/n)$  as a consequence of Bernstein's Inequality (Vershynin, 2018, Corollary 2.8.3) we have that, for any  $t \geq 0$ .

$$P(|\|\mathbf{x}\|_2 - 1| \leq t) \geq 1 - 2 \exp(-cnt^2), \quad (424)$$

for some constant  $c \geq 0$ . Now we have

$$P(\|\mathbf{x}\|_2 \leq g/\sqrt{T}) \begin{cases} \geq 1 - 2 \exp(-cn(g-1)^2) & \text{if } g \geq 1 \\ \leq 2 \exp(-cn(g-1)^2) & \text{otherwise} \end{cases}. \quad (425)$$

We consider the case where  $g \geq 1$ . Then we have that

$$P(\cap_{i=1}^N \mathbf{x}_i \in \mathcal{C}) = P(\cap_{i=1}^N \|\mathbf{x}\|_2 \leq g/\sqrt{T}) \geq 1 - \underbrace{2N \exp(-cn(g-1)^2)}_{=\delta_{\mathcal{C}}}. \quad (426)$$

We have that  $\delta_{\mathcal{C}} = 2N \exp(-cn(g-1)^2)$ .

**Estimating  $B_{\Phi}$ :**

$$B_{\Phi} = \sup_{X \in \mathcal{C}, \{(V_j, \mathbf{z}_j)\} \in \mathcal{F}_{\mathcal{W}}} \left\| \sum_{j=1}^r V_j X \sigma_t(X^T \mathbf{z}_j) \right\|_2, \quad (427)$$

$$\leq R \sup_{X \in \mathcal{C}, \{(V_j, \mathbf{z}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|V_j\|_F \|X \sigma_t(X^T \mathbf{z})\|_2, \quad (428)$$

$$\leq R \sup_{X \in \mathcal{C}, \{(V_j, \mathbf{z}_j)\} \in \mathcal{F}_{\mathcal{W}}} \|V_j\|_F \|X \sigma_t(X^T \mathbf{z})\|_2, \quad (429)$$

$$\leq gRB_V/\sqrt{T}. \quad (430)$$

**Estimating  $B_{\ell}$ :** we have

$$B_{\ell} = \sup_{X \in \mathcal{C}, \{(V_j, \mathbf{z}_j)\} \in \mathcal{F}_{\mathcal{W}}} \left\| \sum_{j=1}^r V_j X \sigma_t(X^T \mathbf{z}_j) - A^* X B^* \right\|, \quad (431)$$

$$= g[RB_V + \|A^*\|_F]/\sqrt{T}. \quad (432)$$

**Estimating  $\tilde{L}_{\Phi}$ :** We have

$$\tilde{L}_{\Phi} = gR\sqrt{B_V^2 + 1/\sqrt{T}}. \quad (433)$$

**Estimating  $\tilde{L}_{\phi}$ :** Similarly we get  $\tilde{L}_{\phi} = g\sqrt{B_V^2 + 1/\sqrt{T}}$  as we have only one slice of factor.

**Estimating  $\epsilon_2$ :** Recall that,

$$\epsilon_2 = \max\{8B_{\ell}\tilde{L}_{\Phi}, 8\tilde{L}_{\Phi}[B_{\ell} + B_{\Phi}L], 32\Omega(f_{\mu}^*)\tilde{L}_{\phi} \max\{B_{\ell}, LB_{\Phi}\}, 4\tilde{L}_{\Phi}B_{\Phi}\}. \quad (434)$$

From all the constants computed earlier, we have that,

$$\epsilon_2 = k_1 g^2 R^2 B_V^2 / T, \quad (435)$$

for some constant  $k_1 \geq 0$ .

Next, we move on estimating  $B(\mathcal{C})$  we need to analyze three terms

**The First Term** is defined via

$$T_1 := \sup_{\{W_j\} \in \mathcal{F}_{\mathcal{W}}} \left| \|f_{\mu}^* \circ \mathcal{P}_{\mathcal{C}} - \Phi_r(\{W_j\}) \circ \mathcal{P}_{\mathcal{C}}\|_{\mu}^2 - \|f_{\mu}^* - \Phi_r(\{W_j\})\|_{\mu}^2 \right|. \quad (436)$$

From Lemma 7 we obtain that, taking  $g \geq 1$  and further simplifying, we get,

$$T_1 \leq c_1 T^2 g e^{-c_2 T g^2} [R^2 \gamma^2 + \|A^*\|_2], \quad (437)$$

**The Second Term** is defined via

$$T_2 := \sup_{\{W_j\} \in \mathcal{F}_W, W' \in \mathcal{F}_\theta} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \phi(W') \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \phi(W') \rangle_\mu \right|. \quad (438)$$

As a consequence of Lemma 7 we have

$$T_2 \leq c_1 T^2 g e^{-c_2 T g^2} [\|A^*\|_2 + R\gamma]. \quad (439)$$

**The Third Term** is defined via

$$T_3 := \sup_{\{W_j\} \in \mathcal{F}_W} \left| \langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, \Phi_r(\{W_j\}) \circ \mathcal{P}_C), \Phi_r(\{W_j\}) \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, \Phi_r(\{W_j\})), \Phi_r(\{W_j\}) \rangle_\mu \right|. \quad (440)$$

As a consequence of Lemma 7 we have

$$T_3 \leq c_1 T^2 g e^{-c_2 T g^2} R\gamma [\|A^*\|_2 + R\gamma], \quad (441)$$

for some positive constant,  $c_1, c_2$ .

Now combining  $T_1, T_2$ , and  $T_3$  from equations (437), (439), (441) we obtain that

$$B(\mathcal{C}) \leq c_1 T^2 g e^{-c_2 T g^2} [\alpha(\|A^*\|_2 + R^2 \gamma^2) + \|A^*\|_2 + R\gamma + R\gamma [\|A^*\|_2 + R\gamma]]. \quad (442)$$

We further upper bound for simplicity as

$$B(\mathcal{C}) \leq 4Rc_1 \sqrt{\frac{\log(T)}{T^5}} g e^{-c_2 g^2} \gamma [\|A^*\|_2 + \gamma]. \quad (443)$$

From Theorem 4 we have that

$$\begin{aligned} \frac{1}{m} |\text{NC}_\mu(\{W_j\}) - \text{NC}_{\mu_N}(\{W_j\})| &\lesssim \frac{1}{2m} \|A^*\|_F \left[ \sup_{\|z\| \leq 1} \frac{1}{N} \left\| \sum_{i=1}^N (Y_i - \hat{Y}_i)^T (X_i \sigma_t(X^T \mathbf{z})) \right\|_F - \lambda \right] \\ &\quad + \frac{2}{m} R c_1 T^2 g e^{-c_2 T g^2} C_V [\|A^*\|_F] [\|A^*\|_2 + \gamma] + C_V^2 [\|A^*\|_F] \\ &\quad \times \sqrt{\frac{R(m+n) \log \left( C_V^2 [\|U^*\|_F^2 + \|V^*\|_F^2]^2 k_1 g^2 R B_V / T \right) \log(N) + \log(1/\delta)}{N}}. \end{aligned} \quad (444)$$

holds true w.p at least  $1 - \delta - 2N \exp(-cn_X(g-1)^2)$ .

Now choose

$$g = 1 + \tilde{O} \left( \frac{1}{T} \log \left( \frac{N}{R(m+n) + \log(1/\delta)} \right) \right). \quad (445)$$

Then we have

$$\frac{2}{m} R c_1 T^2 g e^{-c_2 T g^2} C_V [\|A^*\|_F] [\|A^*\|_2 + \gamma] \lesssim \sqrt{\frac{R(m+n) \log \left( C_V^2 [\|U^*\|_F^2 + \|V^*\|_F^2]^2 k_1 g^2 R B_V / T \right) \log(N) + \log(1/\delta)}{N}}. \quad (446)$$

Therefore we can upper bound the middle term to the right most leaving us behind

$$\begin{aligned} \frac{1}{m} |\mathbf{NC}_\mu(\{W_j\}) - \mathbf{NC}_{\mu_N}(\{W_j\})| &\lesssim \frac{1}{2m} \|A^*\|_F \left[ \frac{1}{N} \sum_{i=1}^N \|Y_i - \hat{Y}_i\|_2 \|X_i\|_2 - \lambda \right] \\ &+ C_V^2 [\|A^*\|_F] \sqrt{\frac{R(m+n) \log(R(m+n)(C_V + B_V)) \log(N) + \log(1/\delta)}{N}}, \end{aligned} \quad (447)$$

holds true w.p at least  $1 - \delta$ .

□

## D GOODS EVENTS

In this section, we provide compute the probabilities of events defined in the proof of Theorem 4. Recall the definition of our function classes:

$$\mathcal{F}_\theta := \{\{W_j\} : \|\Phi_R(\{W_j\})\|_{\text{Lip}} \leq \gamma, \Theta_R(\{W_j\}) \leq \gamma/L_\phi\}; \quad (448)$$

$$\mathcal{F}_\mathcal{W} := \{\{W_j\} : \|\Phi_R(\{W_j\})\|_{\text{Lip}} \leq \gamma, \Theta_R(\{W_j\}) \leq \gamma/L_\phi\}; \quad (449)$$

$$\mathcal{F}_\Phi := \{\Phi_R(\zeta) : \forall \zeta \in \mathcal{F}_\mathcal{W}\} \quad (450)$$

We define the below events:

$$\mathcal{E}_{cvx}(\epsilon) := \{\forall \zeta \in \mathcal{F}_\mathcal{W} : |\mathbf{C}_{\mu_N}(f_\zeta) - \mathbf{C}_\mu(f_\zeta)| \leq \epsilon + B_{nrm}(\mathcal{C})\}; \quad (451)$$

$$\mathcal{E}_{eq}(\epsilon) := \{\forall \zeta \in \mathcal{F}_\mathcal{W} : |\langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_\mu| \leq \epsilon + B_{eq}(\mathcal{C})\}; \quad (452)$$

$$\mathcal{E}_{plr}(\epsilon) := \{\forall \zeta \in \mathcal{F}_\mathcal{W} : |\Omega_{\mu_N}^\circ(\nabla_{\hat{Y}} \ell(g, f_\zeta)) - \Omega_\mu^\circ(\nabla_{\hat{Y}} \ell(g, f_\zeta))| \leq \epsilon + B_{plr}(\mathcal{C})\}; \quad (453)$$

$$\mathcal{E}_{nrm}(\epsilon) := \{\forall \zeta \in \mathcal{F}_\mathcal{W} : \|\mathbf{f}_\mu^* - f_\zeta\|_{\mu_N}^2 - \|\mathbf{f}_\mu^* - f_\zeta\|_\mu^2 \leq \epsilon + B_{nrm}(\mathcal{C})\}. \quad (454)$$

In each of the sections below, we discuss the technical analysis to estimate the probability of the events,  $\mathcal{E}_{cvx}(\epsilon)$ ,  $\mathcal{E}_{eq}(\epsilon)$ ,  $\mathcal{E}_{plr}(\epsilon)$  and  $\mathcal{E}_{nrm}(\epsilon)$ .

### D.1 Concentration of Norms

In this section, we upper bound the probability of the event,  $\mathcal{E}_{nrm}(\epsilon)$  through Lemma 8.

**Lemma 8** (Concentration of Norms). *Consider an  $n_X$ -dimensional sub-Gaussian vector  $X \sim SG(0, (\sigma_X^2/n_X)I)$ , and set of functions  $f_\zeta : \mathbb{R}^{n_X} \rightarrow \mathbb{R}$  as parameterized by  $\zeta \in \mathcal{F}_\mathcal{W}$ . Let  $\mathcal{C}$  be some convex obeying  $P(\bigcap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_C$  for i.i.d samples  $\{X_i\}_{i=1}^N$ . Assume that for any fixed,  $\zeta, \zeta' \in \mathcal{F}_\mathcal{W}$ , and fixed  $Z \in \mathcal{C}$ , we have*

$$\|f_\zeta(Z) - f_{\zeta'}(Z)\| \leq \tilde{L}_\Phi d(\zeta_1, \zeta_2) \text{ and } \|f_\zeta(Z)\| \leq B_\Phi. \quad (455)$$

Denote,

$$B_{nrm}(\mathcal{C}) := \sup_{\zeta \in \mathcal{F}_\mathcal{W}} \|\mathbf{f}_\mu^* \circ \mathcal{P}_\mathcal{C} - f_\zeta \circ \mathcal{P}_\mathcal{C}\|_\mu^2 - \|\mathbf{f}_\mu^* - f_\zeta\|_\mu^2, \quad (456)$$

where  $\mathcal{P}_\mathcal{C}(\cdot)$  denotes the Euclidean projection onto the set  $\mathcal{C}$ . Define,

$$K := 64n_Y\gamma^2\sigma_X^2. \quad (457)$$

Then for any  $\epsilon \in [0, K]$ ,

$$\begin{aligned} & \mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_\mathcal{W}} \|\mathbf{f}_\mu^* - f_\zeta\|_{\mu_N}^2 - \|\mathbf{f}_\mu^* - f_\zeta\|_\mu^2 \geq \epsilon + B_{nrm}(\mathcal{C}) \right) \\ & \leq \delta_C + c \exp \left( \log(C_{\mathcal{F}_\mathcal{W}} \left( \frac{\epsilon}{4\tilde{L}_\Phi B_\Phi} \right)) - N \frac{\epsilon^2}{K^2} \right). \end{aligned} \quad (458)$$

for some positive constant,  $c$  and  $C_{\mathcal{F}_\mathcal{W}}(\nu)$  is the  $\nu$ -net covering number of the set  $\mathcal{F}_\mathcal{W}$ .

*Proof.* If  $X \in \mathbb{R}^{n_X} \sim SG\left(\frac{\sigma_X^2}{n_X} I_{n_X \times n_X}\right)$ , The function map,  $\|\mathbf{f}_\mu^* - f_\zeta\|$  has Lipschitz constant of  $\|\mathbf{f}_\mu^*\|_{\text{Lip}} + \|f_\zeta\|_{\text{Lip}} \leq 2\gamma$ ; as  $f_\zeta, \mathbf{f}_\mu^* \in \mathcal{F}_\Phi$ . Therefore from Theorem 5.1.4 in Vershynin (2018) we have that,  $\mathbf{f}_\mu^*(X) - f_\zeta(X) \sim SG(4\gamma^2\sigma_X^2 I_{n_Y \times n_Y})$ . Thus,  $\|\mathbf{f}_\mu^*(X) - f_\zeta(X)\|^2 \sim SE(4n_Y\gamma^2\sigma_X^2)$

Now, applying the concentration inequality for sub-exponential from Theorem 2.8.1 Vershynin (2018) for a fixed  $\zeta \in \mathcal{F}_\mathcal{W}$ , we have that

$$\mathbb{P}(\|\mathbf{f}_\mu^* - f_\zeta\|_{\mu_N}^2 - \|\mathbf{f}_\mu^* - f_\zeta\|_\mu^2 \geq \epsilon) \leq C \exp \left( -N \min \left\{ \frac{\epsilon^2}{16n_Y^2\gamma^4\sigma_X^4}, \frac{\epsilon}{4n_Y\gamma^2\sigma_X^2} \right\} \right), \quad (459)$$

for some positive constant,  $C \geq 0$ . We use Lemma 14 for applying the concentration bounds. Now set

$$g_\theta = \|f_\mu^* - f_\zeta\|^2. \quad (460)$$

We need to check if the function,  $g$ , is Lipschitz on some metric and convex set  $\mathcal{C} \subseteq \mathbb{R}^{n_X}$ , choose for any  $Z \in \mathcal{C}$ . We have  $P(\bigcap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}$ . Recall that

1.  $\forall \zeta_1, \zeta_2 \in \mathcal{F}_{\mathcal{W}} : \|f_{\zeta_1}(Z) - f_{\zeta_2}(Z)\| \leq \tilde{L}_\Phi d(\zeta_1, \zeta_2)$ , for all  $Z \in \mathcal{C}$ .
2.  $\forall \zeta \in \mathcal{F}_{\mathcal{W}} : \|f_\zeta(Z)\| \leq B_\Phi$ , for all  $Z \in \mathcal{C}$ .
3. For a fixed  $\zeta \in \mathcal{F}_{\mathcal{W}}$ ,

$$|\mathbb{E}[\|f_\mu^*(\mathcal{P}_{\mathcal{C}}(X)) - f_\zeta(\mathcal{P}_{\mathcal{C}}(X))\|^2 - \|f_\mu^*(X) - f_\zeta(X)\|^2]| \leq B_{nrm}(\mathcal{C}). \quad (461)$$

By exploiting the above items we have

$$\begin{aligned} |g_{\theta_1} - g_{\theta_2}| &= \left| \|f_\mu^* - f_{\zeta_1}\|^2 - \|f_\mu^* - f_{\zeta_2}\|^2 \right| = \left| \langle 2f_\mu^* - (f_{\zeta_1} + f_{\zeta_2}), f_{\zeta_1} - f_{\zeta_2} \rangle \right|, \\ &\leq \|2f_\mu^* - (f_{\zeta_1} + f_{\zeta_2})\|^\circ \|f_{\zeta_1} - f_{\zeta_2}\|, \\ &\leq 4\tilde{L}_\Phi B_\Phi d(\zeta_1, \zeta_2). \end{aligned}$$

Then we have that for covering number  $C_{\mathcal{F}_{\mathcal{W}}}(\nu) = \mathcal{N}(\mathcal{F}_{\mathcal{W}}, d(\cdot, \cdot), \nu)$ , and  $K = 4\tilde{L}_\Phi B_\Phi$

$$\begin{aligned} &\mathbb{P}\left(\sup_{\zeta \in \mathcal{F}_{\mathcal{W}}} \left| \|f_\mu^* - f_\zeta\|_{\mu_N}^2 - \|f_\mu^* - f_\zeta\|_\mu^2 \right| \geq \epsilon + B_{nrm} \right) \\ &\leq \delta_{\mathcal{C}} + C \exp\left(\log(C_{\mathcal{F}_{\mathcal{W}}}\left(\frac{\epsilon}{4\tilde{L}_\Phi B_\Phi}\right)) - N \min\left\{\frac{\epsilon^2}{256n_Y^2\gamma^4\sigma_X^4}, \frac{\epsilon}{64n_Y\gamma^2\sigma_X^2}\right\}\right). \end{aligned} \quad (462)$$

We conclude the result by choosing  $\epsilon \in [0, 64n_Y\gamma^2\sigma_X^2]$ .  $\square$

## D.2 Concentration of Convex functions

In this section, we upper bound the probability of the event,  $\mathcal{E}_{cvx}(\epsilon)$  through Lemma 9. In this, we consider strongly and smooth convex function (see assumption 3) through Taylor expansion of the function is always bounded quadratically. Lemma 8 plays an important role in establishing Lemma 9.

**Lemma 9** (Concentration of Convex functions). *Consider an  $n_X$ -dimensional sub-Gaussian vector  $X \sim SG(0, (\sigma_X^2/n_X)I)$ , and set of functions  $f_\zeta : \mathbb{R}^{n_X} \rightarrow \mathbb{R}$  as parameterized by  $\zeta \in \mathcal{F}_{\mathcal{W}}$ . Let  $\mathcal{C}$  be some convex obeying  $P(\bigcap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}$  for i.i.d samples  $\{X_i\}_{i=1}^N$ . Assume that for any fixed,  $\zeta, \zeta' \in \mathcal{F}_{\mathcal{W}}$ , and fixed  $Z \in \mathcal{C}$ , we have*

$$\|f_\zeta(Z) - f_{\zeta'}(Z)\| \leq \tilde{L}_\Phi d(\zeta_1, \zeta_2) \text{ and } \|f_\zeta(Z)\| \leq B_\Phi. \quad (463)$$

Denote

$$B_{nrm}(\mathcal{C}) := \sup_{\zeta \in \mathcal{F}_{\mathcal{W}}} \left| \|f_\mu^* \circ \mathcal{P}_{\mathcal{C}} - f_\zeta \circ \mathcal{P}_{\mathcal{C}}\|_\mu^2 - \|f_\mu^* - f_\zeta\|_\mu^2 \right|. \quad (464)$$

where  $\mathcal{P}_{\mathcal{C}}(\cdot)$  denotes the Euclidean projection onto the set  $\mathcal{C}$ . Define

$$K := n_Y L \left[ (\gamma^2 + \|g\|_{Lip}^2) \sigma_X^2 + \|g\|_{Lip}^2 \sigma_Y^2 \right]. \quad (465)$$

Then for any  $\epsilon \in [0, K]$ ,

$$\begin{aligned} &\mathbb{P}\left(\sup_{\zeta \in \mathcal{F}_{\mathcal{W}}} |C_{\mu_N}(f_\zeta) - C_\mu(f_\zeta)| \geq \epsilon + B_{nrm}(\mathcal{C})\right) \\ &\leq \delta_{\mathcal{C}} + 2 \exp\left(\log\left(C_{\mathcal{F}_{\mathcal{W}}}\left(\frac{\epsilon}{2B_\ell \tilde{L}_\Phi}\right)\right) - cN\left(\frac{\epsilon}{K}\right)^2\right), \end{aligned} \quad (466)$$

for some positive constant,  $c$  and  $C_{\mathcal{F}_{\mathcal{W}}}(\nu)$  is the  $\nu$ -net covering number of the set  $\mathcal{F}_{\mathcal{W}}$ .

*Proof.* Recall the definitions of the convex functions:

$$C_{\mu_N}(f) := \ell(g, f)_{\mu_N} + \lambda \Omega(f), \text{ and } C_{\mu}(f) := \ell(g, f)_{\mu} + \lambda \Omega(f). \quad (467)$$

The difference between these two terms is

$$|C_{\mu_N}(f) - C_{\mu}(f)| = |\ell(g, f)_{\mu_N} - \ell(g, f)_{\mu}|. \quad (468)$$

From assumption 3,  $\ell(\cdot, \cdot)$  is second-order differentiable in the second argument. By 2nd-order Taylor's theorem, we have

$$\ell(Y, \hat{Y}) = \ell(Y, \hat{Y}_0) + \langle \nabla_{\hat{Y}} \ell(Y, \hat{Y}_0), \hat{Y} - \hat{Y}_0 \rangle + \langle \int_0^1 t \nabla_{\hat{Y}}^2 \ell(Y, \hat{Y}_0 + t(\hat{Y} - \hat{Y}_0)) dt, (\hat{Y} - \hat{Y}_0)(\hat{Y} - \hat{Y}_0)^T \rangle. \quad (469)$$

Now choose  $Y = \hat{Y}_0 = g(X(\omega), E(\omega))$ , and  $\hat{Y} = f_{\zeta}(X(\omega))$ . As  $\ell(Y, Y) = 0$ , and  $\nabla_{\hat{Y}} \ell(Y, Y) = \mathbf{0}$ . Plugging these parameters in the Taylor expansion we have that (ignoring the inputs,  $(X(\omega), E(\omega))$  for simplicity),

$$\ell(g, f) = \langle \int_0^1 t \nabla_{\hat{Y}}^2 \ell(g, g + t(f_{\zeta} - g)) dt, (f_{\zeta} - g)(f_{\zeta} - g)^T \rangle. \quad (470)$$

Now we apply expectation over the measure  $\mu_N$ , and  $\mu$  respectively on the above equality. Then we have that

$$\ell(g, f_{\zeta})_{\mu} = \langle \int_0^1 t \nabla_{\hat{Y}}^2 \ell(g, g + t(f_{\zeta} - g)) dt, (f_{\zeta} - g)(f_{\zeta} - g)^T \rangle_{\mu}; \quad (471)$$

$$\ell(g, f_{\zeta})_{\mu_N} = \langle \int_0^1 t \nabla_{\hat{Y}}^2 \ell(g, g + t(f_{\zeta} - g)) dt, (f_{\zeta} - g)(f_{\zeta} - g)^T \rangle_{\mu_N}. \quad (472)$$

Since  $f_{\zeta}$  and  $g$  are Lipschitz functions and the inputs are sub-Gaussian, we have that  $f_{\zeta} - g$  is a sub-Gaussian vector. As a consequence of Lemma 2.7.6 from Vershynin (2018) we obtain that  $(f_{\zeta} - g)(f_{\zeta} - g)^T$  follows a sub-exponential distribution, whose concentration is well-studied.

As a consequence of assumption 3 the hessian is bounded, i.e.,  $\alpha I \preceq \nabla_{\hat{Y}}^2 \ell(\cdot, \cdot) \preceq LI$ . We can argue that the product of a bounded RV and sub-exponential RV is sub-exponential. Recall Item (iii) from Proposition 2.7.1 of Vershynin (2018). The random variable  $Z$  is sub-exponential iff

$$\mathbb{E}_Z [e^{\lambda|Z|}] \leq e^{\lambda K}; \forall \lambda \in [0, 1/K], \quad (473)$$

for some positive constant,  $K \geq 0$ .

We now verify if  $\langle H(\mathbf{x}, \mathbf{z}), \mathbf{x}\mathbf{x}^T \rangle$  is sub-exponential. Given that  $\mathbf{x} \sim SG(\sigma_X^2/n_x I_{n_x \times n_x})$ ,  $\mathbf{z}$  is a R.V. Suppose  $A \preceq H(\mathbf{x}, \mathbf{z}) \preceq B$  a.s. Then we have that

$$\mathbb{E}_{\mathbf{x}, \mathbf{z}} [e^{\lambda |\langle H(\mathbf{x}, \mathbf{z}), \mathbf{x}\mathbf{x}^T \rangle|}] \leq \mathbb{E}_{\mathbf{x}, \mathbf{z}} [e^{\lambda \|H(\mathbf{x}, \mathbf{z})\|_2 \|\mathbf{x}\mathbf{x}^T\|_2}], \quad (474)$$

$$\leq \mathbb{E}_{\mathbf{x}, \mathbf{z}} [e^{\lambda \max\{\rho(A), \rho(B)\} \|\mathbf{x}\|_2^2}], \quad (475)$$

$$(476)$$

where,  $\rho(A)$  is the spectral radius of the matrix,  $A$ . Since,  $\mathbf{x} \sim SG(\sigma_X^2/n_x I_{n_x \times n_x})$ , we have  $\|\mathbf{x}\|^2 \sim SE(\sigma_X^2)$ . Then,

$$\mathbb{E}_{\mathbf{x}, \mathbf{z}} [e^{\lambda |\langle H(\mathbf{x}, \mathbf{z}), \mathbf{x}\mathbf{x}^T \rangle|}] \leq e^{\lambda \max\{\rho(A), \rho(B)\} \sigma_X^2},$$

implies that  $\langle H(\mathbf{x}, \mathbf{z}), \mathbf{x}\mathbf{x}^T \rangle \sim SE(\max\{\rho(A), \rho(B)\} \sigma_X^2)$ . From this analysis we have

$$\begin{aligned} & \langle \int_0^1 t \nabla_{\hat{Y}}^2 \ell(g, g + t(f_{\zeta} - g)) dt, \underbrace{(f_{\zeta} - g)(f_{\zeta} - g)^T}_{\sim SE(\left[ (\|f_{\zeta}\|_{\text{Lip}}^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right] I_{n_y \times n_y}} \rangle \\ & \sim SE\left(\left[ (\|f_{\zeta}\|_{\text{Lip}}^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right] I_{n_y \times n_y}\right) \end{aligned}$$



$$\sim SE \left( n_Y \frac{L}{2} \left[ (\|f_\zeta\|_{\text{Lip}}^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right] \right). \quad (477)$$

For convex functions, we know that  $0 \leq \alpha \leq L$ . As a consequence, we have  $\frac{\alpha}{2} I \preceq \int_0^1 t \nabla_Y^2 \ell(g, g+t(f_\zeta-g)) dt \preceq \frac{L}{2} I$ . Now we apply sub-exponential concentration for a fixed  $\zeta \in \mathcal{F}_W$ , yielding

$$\begin{aligned} & \mathbb{P}(|C_{\mu_N}(f_\zeta) - C_\mu(f_\zeta)| \geq \epsilon) \\ & \leq 2 \exp \left( -cN \min \left\{ \left( \frac{2\epsilon}{n_Y L \left[ (\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right]} \right)^2, \right. \right. \\ & \quad \left. \left. \frac{2\epsilon}{n_Y L \left[ (\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right]} \right\} \right), \end{aligned} \quad (478)$$

for some positive constant,  $c \geq 0$ .

Next, we move on to obtain a uniform concentration for all  $\zeta \in \mathcal{F}_W$ . Now we apply covering argument from Lemma 14, and set

$$g_\theta = \ell(g, f_\zeta). \quad (479)$$

We need to check if the function,  $g$ , is Lipschitz on some metric and convex set  $\mathcal{C} \subseteq \mathbb{R}^{n_X}$ , choose for any  $Z \in \mathcal{C}$ . We have  $P(\bigcap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}$ . Recall that

1.  $\forall \zeta_1, \zeta_2 \in \mathcal{F}_W : \|f_{\zeta_1}(Z) - f_{\zeta_2}(Z)\| \leq \tilde{L}_\Phi d(\zeta_1, \zeta_2)$ , for all  $Z \in \mathcal{C}$ .
2.  $\forall \zeta \in \mathcal{F}_W : \|\nabla_Y \ell(g(Z), f_\zeta(Z))\| \leq B_\ell$ , for all  $Z \in \mathcal{C}$ .
3. For a fixed  $\zeta \in \mathcal{F}_W$ ,

$$|\mathbb{E} [\|f_\mu^*(\mathcal{P}_{\mathcal{C}}(X)) - f_\zeta(\mathcal{P}_{\mathcal{C}}(X))\|^2 - \|f_\mu^*(X) - f_\zeta(X)\|^2]| \leq B_{nrm}(\mathcal{C}). \quad (480)$$

From Taylor expansion we have that,

$$\begin{aligned} |g_{\theta_1} - g_{\theta_2}| &= |\langle \int_t \nabla_Y \ell(g, f_{\zeta_1} + t(f_{\zeta_2} - f_{\zeta_1})) dt, f_{\zeta_1} - f_{\zeta_2} \rangle|, \\ &\leq \left\| \int_t \nabla_Y \ell(g, f_{\zeta_1} + t(f_{\zeta_2} - f_{\zeta_1})) dt \right\| \|f_{\zeta_2} - f_{\zeta_1}\|, \\ &\leq B_\ell \|f_{\zeta_2} - f_{\zeta_1}\|, \\ &\leq B_\ell \tilde{L}_\Phi d(\zeta_1, \zeta_2). \end{aligned}$$

From Lemma 14 we have

$$\mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_W} |C_{\mu_N}(f_\zeta) - C_\mu(f_\zeta)| \geq \epsilon \right) \quad (481)$$

$$\begin{aligned} & \leq \delta_{\mathcal{C}} + 2 \exp \left( \log \left( C_{\mathcal{F}_W} \left( \frac{\epsilon}{2B_\ell \tilde{L}_\Phi} \right) \right) - cN \min \left\{ \left( \frac{\epsilon}{n_Y L \left[ (\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right]} \right)^2, \right. \right. \\ & \quad \left. \left. \frac{\epsilon}{n_Y L \left[ (\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right]} \right\} \right), \end{aligned} \quad (482)$$

for some positive constant,  $c$ . Now restrict  $\epsilon \in \left[ 0, n_Y L \left[ (\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{Y|X}^2 \right] \right]$ . This completes our proof.  $\square$

### D.3 Concentration of Equilibria

In this section, we upper bound the probability of the event,  $\mathcal{E}_{eq}(\epsilon)$  through Lemma 10. We first present the concentration of bi-Lipschitz functions for sub-Gaussian inputs.

**Proposition 4.** *Let  $X \sim SG(\frac{\sigma_X^2}{n_x} I_{n_x \times n_x})$  and  $Y|X \sim SG(\frac{\sigma_{Y|X}^2}{n_y} I_{n_y \times n_y})$ , where  $\sigma_{Y|X}$  is independent of  $X = x$ , then for any Lipschitz function,  $\phi : \mathcal{Z} \rightarrow \mathbb{R}$  and a function,  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Z}$ , satisfies the following. Then  $\phi(f(X, Y)) \sim SG\left(4\|\phi\|_{Lip}^2\|f\|_{Lip}^2\left[\sigma_X^2 + \sigma_{Y|X}^2\right]\right)$ .*

$$\|f(X_2, Y_2) - f(X_1, Y_1)\| \leq \|f\|_{Lip} [\|X_2 - X_1\| + \|Y_2 - Y_1\|]. \quad (483)$$

*Proof.* Let us compute the moments of the random variable  $\phi(f(X, Y))$  for any Lipschitz function,  $\phi : \mathcal{Z} \rightarrow \mathbb{R}$ .

By symmetrization, we have

$$\begin{aligned} & \mathbb{E}_{X,Y} [\exp(\lambda[\phi(f(X, Y)) - \mathbb{E}_{X,Y}[\phi(f(X, Y))])]] \\ &= \mathbb{E}_{X,Y} [\exp(\lambda[\phi(f(X, Y)) - \mathbb{E}_{X',Y'}[\phi(f(X', Y'))])]] \end{aligned} \quad (484)$$

$$= \mathbb{E}_{X,Y} [\exp(\lambda \mathbb{E}_{X',Y'} [\phi(f(X, Y)) - \phi(f(X', Y'))])]. \quad (485)$$

By Jensen's inequality we have

$$\leq \mathbb{E}_{X,Y,X',Y'} [\exp(\lambda[\phi(f(X, Y)) - \phi(f(X', Y'))])], \quad (486)$$

By Lipschitz continuity we have

$$\leq \mathbb{E}_{X,Y,X',Y'} [\exp(\|\phi\|_{Lip} \lambda [\|f(X, Y) - f(X', Y')\|])]. \quad (487)$$

By construction, we have.

$$\leq \mathbb{E}_{X,Y,X',Y'} [\exp(\|\phi\|_{Lip} \|f\|_{Lip} \lambda [\|X - X'\| + \|Y - Y'\|])]. \quad (488)$$

By Cauchy-Schwartz's inequality we have

$$\leq \mathbb{E}_{X,Y,X',Y'} [\exp(\|\phi\|_{Lip} \|f\|_{Lip} \lambda [\|X\| + \|X'\| + \|Y\| + \|Y'\|])]. \quad (489)$$

As the symmetrized random variables are independent,

$$\leq \mathbb{E}_{X,Y} [\exp(2\|\phi\|_{Lip} \|f\|_{Lip} \lambda [\|X\| + \|Y'\|])]. \quad (490)$$

Now perform conditional expectation,

$$\leq \mathbb{E}_X \mathbb{E}_{Y|X} [\exp(2\|\phi\|_{Lip} \|f\|_{Lip} \lambda [\|X\| + \|Y\|])], \quad (491)$$

$$\leq \mathbb{E}_X \exp(2\|\phi\|_{Lip} \|f\|_{Lip} \lambda [\|X\|]) \mathbb{E}_{Y|X} [\exp(2\|\phi\|_{Lip} \|f\|_{Lip} \lambda [\|Y\|])], \quad (492)$$

$$\leq \exp\left(\frac{\lambda^2}{2} 4\|\phi\|_{Lip}^2 \|f\|_{Lip}^2 [\sigma_X^2 + \sigma_{Y|X}^2] + \lambda [\mathbb{E}_X[\|X\|] + \mathbb{E}_{Y|X}[\|Y\|]]\right), \quad (493)$$

$$\leq K \exp\left(\frac{\lambda^2}{2} 4\|\phi\|_{Lip}^2 \|f\|_{Lip}^2 [\sigma_X^2 + \sigma_{Y|X}^2]\right). \quad (494)$$

for some constant,  $K \geq 0$ .

This implies that,  $\phi(f(X, Y)) \sim SG\left(4\|\phi\|_{Lip}^2\|f\|_{Lip}^2\left[\sigma_X^2 + \sigma_{Y|X}^2\right]\right)$ .  $\square$

With the above result, we now state and prove the probability of the event,  $\mathcal{E}_{eq}(\epsilon)$ .

**Lemma 10** (Concentration of Equilibria). *Consider an  $n_X$ -dimensional sub-Gaussian vector  $X \sim SG(0, (\sigma_X^2/n_X)I)$ , and set of functions  $f_\zeta : \mathbb{R}^{n_X} \rightarrow \mathbb{R}$  as parameterized by  $\zeta \in \mathcal{F}_W$ . Let  $\mathcal{C}$  be some convex obeying  $P(\bigcap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_C$  for i.i.d samples  $\{X_i\}_{i=1}^N$ . Assume that for any fixed,  $\zeta_1, \zeta_2 \in \mathcal{F}_W$ , and fixed  $Z \in \mathcal{C}$ , we have*

$$\|\nabla_{\hat{Y}} \ell(g(Z), f_\zeta(Z))\| \leq B_\ell, \|f_\zeta(Z)\| \leq B_\Phi, \text{ and} \quad (495)$$

$$\|f_\zeta(Z) - f_{\zeta'}(Z)\| \leq \tilde{L}_\Phi d(\zeta, \zeta'). \quad (496)$$

In addition, we have that,

$$\sup_{\zeta \in \mathcal{F}_W} |\mathbb{E} [\langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, f_\zeta \circ \mathcal{P}_C), f_{\zeta'} \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_\mu ]| = B_{eq}(C). \quad (497)$$

Define

$$K := 4n_y \gamma \|\nabla_{\hat{Y}} \ell\|_{Lip} \sigma_X \sqrt{(\gamma^2 + \|g\|_{Lip}^2) \sigma_X^2 + \|g\|_{Lip}^2 \sigma_{E|X}^2}. \quad (498)$$

Then for any  $\epsilon \in [0, K]$ ,

$$\begin{aligned} \mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_W} |\langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_\mu | \geq \epsilon + B_{eq}(C) \right) \leq \\ \delta_C + c \exp \left( \log \left( C_{\mathcal{F}_W} \left( \frac{\epsilon}{2\tilde{L}_\Phi [B_\ell + B_\Phi L]} \right) \right) - N \frac{\epsilon^2}{K^2} \right), \end{aligned} \quad (499)$$

for some positive constant,  $c$  and  $C_{\mathcal{F}_W}(\nu)$  is the  $\nu$ -net covering number of the set  $\mathcal{F}_W$ .

*Proof.* From Assumptions 1-5, we have that for any  $g_1, f_1, g_2, f_2 \in L^2(\mu)$

$$\|\nabla_{\hat{Y}} \ell(g_2(X(\omega), E(\omega)), f_2(X(\omega))) - \nabla_{\hat{Y}} \ell(g_1(X(\omega), E(\omega)), f_1(X(\omega)))\| \quad (500)$$

$$\leq \|\nabla_{\hat{Y}} \ell\|_{Lip} [\|g_2(X(\omega), E(\omega)) - g_1(X(\omega), E(\omega))\| + \|f_2(X(\omega)) - f_1(X(\omega))\|]. \quad (501)$$

Since  $X(\omega)$  and  $E(\omega)$  are Lipschitz concentrated R.Vs, it holds that

$$g(X, E)|E \sim SG(\|g\|_{Lip}^2 \sigma_X^2 I_{n_y \times n_y}), g(X, E)|X \sim SG(\|g\|_{Lip}^2 \sigma_{E|X}^2 I_{n_y \times n_y}), \quad (502)$$

$$\text{and } f_\zeta(X) \sim SG(\|f_\zeta\|_{Lip}^2 \sigma_X^2 I_{n_y \times n_y}). \quad (503)$$

From Proposition 4 we have

$$\begin{aligned} \nabla_{\hat{Y}} \ell(g(X(\omega), E(\omega)), f_\zeta(X(\omega))) \sim \\ SG \left( 4\|\nabla_{\hat{Y}} \ell\|_{Lip}^2 \left[ (\|f_\zeta\|_{Lip}^2 + \|g\|_{Lip}^2) \sigma_X^2 + \|g\|_{Lip}^2 \sigma_{E|X}^2 \right] I_{n_y} \right). \end{aligned} \quad (504)$$

Now we have the inner product between two sub-Gaussian random variables from Proposition 7 we have that the result is sub-exponential, i.e.,

$$\begin{aligned} \langle \underbrace{\nabla_{\hat{Y}} \ell(g(X(\omega), E(\omega)), f_\zeta(X(\omega)))}_{\sim SG(4\|\nabla_{\hat{Y}} \ell\|_{Lip}^2 [(\|f_\zeta\|_{Lip}^2 + \|g\|_{Lip}^2) \sigma_X^2 + \|g\|_{Lip}^2 \sigma_{E|X}^2] I_{n_y \times n_y})}, \underbrace{f_\zeta(X(\omega))}_{\sim SG(\|f_\zeta\|_{Lip}^2 \sigma_X^2 I_{n_y})} \rangle \\ \sim SE \left( 2n_y \|\nabla_{\hat{Y}} \ell\|_{Lip} \|f_\zeta\|_{Lip} \sigma_X \sqrt{(\|f_\zeta\|_{Lip}^2 + \|g\|_{Lip}^2) \sigma_X^2 + \|g\|_{Lip}^2 \sigma_{E|X}^2} \right). \end{aligned} \quad (505)$$

The class of functions,  $f_\zeta$  for  $\zeta \in \mathcal{F}_W$ , has bounded Lipschitz constant  $\gamma$ . As a consequence of the sub-exponential concentration bound from Theorem 2.8.1 in Vershynin (2018), we have that for a fixed  $\zeta \in \mathcal{F}_W$ ,

$$\mathbb{P} (|\langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_\mu | \geq \epsilon) \leq C \exp \left( -N \min \left\{ \frac{\epsilon^2}{K^2}, \frac{\epsilon}{K} \right\} \right). \quad (506)$$

where  $K := 2n_Y \gamma \|\nabla_{\hat{Y}} \ell\|_{\text{Lip}} \sigma_X \sqrt{(\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{E|X}^2}$  and some positive constant,  $C$ .

Now we move on to providing a uniform concentration in the inequality (506). We will apply uniform concentration result from Lemma 14, for this set:

$$g_\theta = \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle. \quad (507)$$

Recall the below items:

1. For a fixed  $Z \in \mathcal{C}$  we have  $\forall \zeta \in \mathcal{F}_W : \|\nabla_{\hat{Y}} \ell(g(Z), f_\zeta(Z))\| \leq B_\ell$ .
2. For a fixed  $Z \in \mathcal{C}$  we have  $\forall \zeta \in \mathcal{F}_W : \|f_\zeta(Z)\| \leq B_\Phi$ .
3. For a fixed  $Z \in \mathcal{C}$  we have  $\forall \zeta, \zeta' \in \mathcal{F}_W : \|f_\zeta(Z) - f_{\zeta'}(Z)\| \leq \tilde{L}_\Phi d(\zeta, \zeta')$ .
4. For a any  $\hat{Y}_1, \hat{Y}_2 \in \mathbb{R}^{n_Y}$  we have  $\|\nabla_{\hat{Y}} \ell(Y, \hat{Y}_1) - \nabla_{\hat{Y}} \ell(Y, \hat{Y}_2)\| \leq L \|\hat{Y}_1 - \hat{Y}_2\|$ .
5. For a fixed  $\zeta \in \mathcal{F}_W$ ,

$$\sup_{\zeta \in \mathcal{F}_W} \left| \mathbb{E} [\langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, f_\zeta \circ \mathcal{P}_C), f_\zeta \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_\mu] \right| = B_{eq}(\mathcal{C}). \quad (508)$$

Now we check the Lipschitz continuity of the function  $g_\theta$ :

$$\begin{aligned} |g_{\theta_1} - g_{\theta_2}| &= |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta_1} \rangle - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_2}), f_{\zeta_2} \rangle|, \\ &= |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta_1} - f_{\zeta_2} \rangle - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_2}) - \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta_2} \rangle|, \\ &\leq |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta_1} - f_{\zeta_2} \rangle| + |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_2}) - \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta_2} \rangle|, \\ &\leq \|\nabla_{\hat{Y}} \ell(g, f_{\zeta_1})\| \|f_{\zeta_1} - f_{\zeta_2}\| + \|\nabla_{\hat{Y}} \ell(g, f_{\zeta_2}) - \nabla_{\hat{Y}} \ell(g, f_{\zeta_1})\| \|f_{\zeta_2}\|, \\ &\leq B_\ell \|f_{\zeta_1} - f_{\zeta_2}\| + B_\Phi L \|f_{\zeta_1} - f_{\zeta_2}\|, \\ &\leq \tilde{L}_\Phi [B_\ell + B_\Phi L] d(\zeta_1, \zeta_2). \end{aligned}$$

Then from Lemma 14 we have that

$$\mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_W} |\langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_\zeta \rangle_\mu| \geq \epsilon + B_{eq}(\mathcal{C}) \right) \quad (509)$$

$$\leq \delta_C + C \exp \left( \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{2\tilde{L}_\Phi [B_\ell + B_\Phi L]} \right) \right) - N \min \left\{ \frac{\epsilon^2}{4K^2}, \frac{\epsilon}{2K} \right\} \right). \quad (510)$$

□

#### D.4 Concentration of Polar

In this section, we compute the probability of the occurrence of the event,  $\mathcal{E}_{plr}(\epsilon)$  through Lemma 11. The analysis of  $\mathcal{E}_{plr}(\epsilon)$  resembles to that of  $\mathcal{E}_{eq}(\epsilon)$  following similar arguments.

**Lemma 11** (Concentration of Polar). *Consider an  $n_X$ -dimensional sub-Gaussian vector  $X \sim SG(0, (\sigma_X^2/n_X)I)$ , and set of functions  $f_\zeta : \mathbb{R}^{n_X} \rightarrow \mathbb{R}$  as parameterized by  $\zeta \in \mathcal{F}_W$ . Let  $\mathcal{C}$  be some convex obeying  $P(\bigcap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_C$  for i.i.d samples  $\{X_i\}_{i=1}^N$ . Assume that for any fixed,  $\zeta_1, \zeta_2 \in \mathcal{F}_W$ ,  $\zeta'_1, \zeta'_2 \in \mathcal{F}_\theta$ , and fixed  $Z \in \mathcal{C}$ , we have*

$$\|\nabla_{\hat{Y}} \ell(g(Z), f_\zeta(Z))\| \leq B_\ell, \|f_\zeta(Z)\| \leq B_\Phi, \quad (511)$$

$$\|f_{\zeta_1}(Z) - f_{\zeta_2}(Z)\| \leq \tilde{L}_\Phi d(\zeta_1, \zeta_2), \text{ and } \|f_{\zeta'_1}(Z) - f_{\zeta'_2}(Z)\| \leq \tilde{L}_\Phi d(\zeta'_1, \zeta'_2). \quad (512)$$

In addition, we have that,

$$\sup_{\zeta \in \mathcal{F}_W, \zeta' \in \mathcal{F}_\theta} \left| \mathbb{E} [\langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_C, f_\zeta \circ \mathcal{P}_C), f_{\zeta'} \circ \mathcal{P}_C \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_\mu] \right| = B_{plr}(\mathcal{C}). \quad (513)$$

Define

$$K := 4n_Y \|\nabla_{\hat{Y}} \ell\|_{\text{Lip}} L_\phi \sigma_X \sqrt{(\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{E|X}^2}. \quad (514)$$

Then for any  $\epsilon \in [0, K]$ ,

$$\begin{aligned} & \mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_{\mathcal{W}}} : |\Omega_{\mu_N}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) - \Omega_{\mu}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta}))| \geq \epsilon + B_{plr}(\mathcal{C}) \right) \\ & \leq \delta_{\mathcal{C}} + c \exp \left( \log \left( \mathcal{C}_{\mathcal{F}_{\mathcal{W}}} \left( \frac{\epsilon}{8 \max\{\tilde{L}_{\phi} B_{\ell}, L \tilde{L}_{\Phi} B_{\Phi}\}} \right) \right) \right. \\ & \quad \left. + \log \left( \mathcal{C}_{\mathcal{F}_{\theta}} \left( \frac{\epsilon}{8 \max\{\tilde{L}_{\phi} B_{\ell}, L \tilde{L}_{\Phi} B_{\Phi}\}} \right) \right) - N \frac{\epsilon^2}{K^2} \right), \end{aligned} \quad (515)$$

for some positive constant,  $c$  and  $\mathcal{C}_{\mathcal{F}_{\mathcal{W}}}(\nu)$  (and  $\mathcal{C}_{\mathcal{F}_{\theta}}(\nu)$ ) is the  $\nu$ -net covering number of the set  $\mathcal{F}_{\mathcal{W}}$  (and  $\mathcal{F}_{\theta}$ ).

*Proof.* Recall the definition of the polar in Equation 6:

$$\Omega_{\mu_N}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) := \sup_{\zeta' \in \mathcal{F}_{\theta}} \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu_N}, \quad (516)$$

$$\Omega_{\mu}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) := \sup_{\zeta' \in \mathcal{F}_{\theta}} \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu}. \quad (517)$$

Now, by taking the difference between the above two polars, we have

$$\begin{aligned} & \left| \Omega_{\mu_N}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) - \Omega_{\mu}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) \right| \\ & = \left| \sup_{\zeta' \in \mathcal{F}_{\theta}} \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu_N} - \sup_{\zeta' \in \mathcal{F}_{\theta}} \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu} \right|. \end{aligned} \quad (518)$$

Denote,  $\zeta_{\mu}^{t*} = \arg \sup_{\zeta' \in \mathcal{F}_{\theta}} \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu}$  and  $\zeta_{\mu_N}^{t*} = \arg \sup_{\zeta' \in \mathcal{F}_{\theta}} \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu_N}$ , then, by definition we have that

$$\begin{aligned} & -\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu}^{t*}} \rangle_{\mu_N} + \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu}^{t*}} \rangle_{\mu} \leq \Omega_{\mu_N}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) - \Omega_{\mu}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) \leq \\ & \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu_N}^{t*}} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu}^{t*}} \rangle_{\mu}. \end{aligned} \quad (519)$$

Applying modulus on both sides we obtain

$$\begin{aligned} & \left| \Omega_{\mu_N}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) - \Omega_{\mu}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) \right| \\ & \leq \max \left\{ \left| \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu_N}^{t*}} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu}^{t*}} \rangle_{\mu} \right|, \right. \\ & \quad \left| \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu_N}^{t*}} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta_{\mu}^{t*}} \rangle_{\mu} \right\} \\ & \leq \sup_{\zeta' \in \mathcal{F}_{\theta}} \left| \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu} \right|. \end{aligned} \quad (520)$$

Now, we have to compute a lower bound on

$$\mathbb{P} \left( \sup_{\zeta' \in \mathcal{F}_{\theta}} \left| \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta}), f_{\zeta'} \rangle_{\mu} \right| \leq \epsilon \right). \quad (521)$$

The computation of Equation (521) is similar to that of Lemma 10. We can re-write the concentration of polars and apply the monotonicity of probability in inequality (520) by doing this we have

$$\mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_{\mathcal{W}}} : |\Omega_{\mu_N}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta})) - \Omega_{\mu}^{\circ}(\nabla_{\hat{Y}} \ell(g, f_{\zeta}))| \geq \epsilon \right)$$

$$\leq \mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_\Phi, \zeta' \in \mathcal{F}_\theta} : |\langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_\mu | \leq \epsilon \right). \quad (522)$$

As the data follow the sub-Gaussian distribution from Proposition ?? and Proposition 7, we have

$$\begin{aligned} & \left\langle \underbrace{\nabla_{\hat{Y}} \ell(g, f_\zeta)}_{\sim SG(4\|\nabla_{\hat{Y}} \ell\|_{\text{Lip}}^2 \left[ (\|f_\zeta\|_{\text{Lip}}^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{E|X}^2 \right] I_{n_Y \times n_Y}}}, \underbrace{f_{\zeta'}}_{\sim SG(\|f_{\zeta'}\|_{\text{Lip}}^2 \sigma_X^2 I_{n_Y \times n_Y})} \right\rangle_{\mu_N} \\ & \sim SE \left( 2n_Y \|\nabla_{\hat{Y}} \ell\|_{\text{Lip}} \|f_{\zeta'}\|_{\text{Lip}} \sigma_X \sqrt{(\|f_\zeta\|_{\text{Lip}}^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{E|X}^2} \right). \end{aligned} \quad (523)$$

From Assumption 6 the class,  $\mathcal{F}_\Phi$  has Lipschitz constant at most,  $\gamma$ . From the assumption 4  $\mathcal{F}_\theta$  has a Lipschitz constant at most  $\gamma_\theta$ . Therefore, the inner product described above is concentrated as a consequence of Theorem 2.8.1 from Vershynin (2018). Now for a fixed  $\zeta \in \mathcal{F}_\mathcal{W}, \zeta' \in \mathcal{F}_\theta$ , we have that

$$\mathbb{P} \left( |\langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_\mu | \leq \epsilon \right) \leq C \exp \left( -N \min \left\{ \frac{\epsilon^2}{K^2}, \frac{\epsilon}{K} \right\} \right). \quad (524)$$

where,  $K = 2n_Y \|\nabla_{\hat{Y}} \ell\|_{\text{Lip}} L_\phi \sigma_X \sqrt{(\gamma^2 + \|g\|_{\text{Lip}}^2) \sigma_X^2 + \|g\|_{\text{Lip}}^2 \sigma_{E|X}^2}$ .

Now we utilize Lemma 14 to have this concentration uniformly for all,  $\zeta \in \mathcal{F}_\mathcal{W}, \zeta' \in \mathcal{F}_\theta$ . Set

$$g_\theta = \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle. \quad (525)$$

Recall the below items:

1. For a fixed  $Z \in \mathcal{C}$  we have  $\forall \zeta \in \mathcal{F}_\mathcal{W} : \|\nabla_{\hat{Y}} \ell(g(Z), f_\zeta(Z))\| \leq B_\ell$ .
2. For a fixed  $Z \in \mathcal{C}$  we have  $\forall \zeta \in \mathcal{F}_\mathcal{W} : \|f_\zeta(Z)\| \leq B_\Phi$ .
3. For a fixed  $Z \in \mathcal{C}$  we have  $\forall \zeta, \zeta' \in \mathcal{F}_\mathcal{W} : \|f_\zeta(Z) - f_{\zeta'}(Z)\| \leq \tilde{L}_\Phi d(\zeta, \zeta')$ .
4. For a fixed  $Z \in \mathcal{C}$  we have  $\forall \zeta, \zeta' \in \mathcal{F}_\theta : \|f_\zeta(Z) - f_{\zeta'}(Z)\| \leq \tilde{L}_\phi d(\zeta, \zeta')$ .
5. For a any  $\hat{Y}_1, \hat{Y}_2 \in \mathbb{R}^{n_Y}$  we have  $\|\nabla_{\hat{Y}} \ell(Y, \hat{Y}_1) - \nabla_{\hat{Y}} \ell(Y, \hat{Y}_2)\| \leq L \|\hat{Y}_1 - \hat{Y}_2\|$ .
6. For a fixed  $\zeta \in \mathcal{F}_\mathcal{W}, \zeta' \in \mathcal{F}_\theta$ ,

$$\sup_{\zeta \in \mathcal{F}_\mathcal{W}, \zeta' \in \mathcal{F}_\theta} \left| \mathbb{E} [\langle \nabla_{\hat{Y}} \ell(g \circ \mathcal{P}_\mathcal{C}, f_\zeta \circ \mathcal{P}_\mathcal{C}), f_{\zeta'} \circ \mathcal{P}_\mathcal{C} \rangle_\mu - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_\mu] \right| = B_{plr}(\mathcal{C}). \quad (526)$$

Now we check the Lipschitzness of  $g$ :

$$\begin{aligned} |g_{\theta_1} - g_{\theta_2}| &= |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta'_1} \rangle - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_2}), f_{\zeta'_2} \rangle| \\ &= |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta'_1} - f_{\zeta'_2} \rangle - \langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_2}) - \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta'_2} \rangle|, \\ &\leq |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta'_1} - f_{\zeta'_2} \rangle| + |\langle \nabla_{\hat{Y}} \ell(g, f_{\zeta_2}) - \nabla_{\hat{Y}} \ell(g, f_{\zeta_1}), f_{\zeta'_2} \rangle|, \\ &\leq \|\nabla_{\hat{Y}} \ell(g, f_{\zeta_1})\| \|f_{\zeta'_1} - f_{\zeta'_2}\| + \|\nabla_{\hat{Y}} \ell(g, f_{\zeta_2}) - \nabla_{\hat{Y}} \ell(g, f_{\zeta_1})\| \|f_{\zeta'_2}\|, \\ &\leq B_\ell \tilde{L}_\phi d(\zeta'_1, \zeta'_2) + B_\Phi L \tilde{L}_\Phi d(\zeta_1, \zeta_2), \\ &\leq 2 \max\{\tilde{L}_\phi B_\ell, L \tilde{L}_\Phi B_\Phi\} \max\{d(\zeta'_1, \zeta'_2), d(\zeta_1, \zeta_2)\}. \end{aligned}$$

Now, we have a product of two metric spaces whose metric is a maximum of individual metrics, therefore simply we can upper bound the covering number by product of these two metric spaces, i.e.,

$$\mathcal{N}(\mathcal{F}_\mathcal{W} \times \mathcal{F}_\theta, \|\cdot\|_{\infty, d(\cdot, \cdot)}, \nu) \leq \mathcal{N}(\mathcal{F}_\mathcal{W}, d(\cdot, \cdot), \nu) \mathcal{N}(\mathcal{F}_\theta, d(\cdot, \cdot), \nu). \quad (527)$$

From Lemma 14 we have that

$$\mathbb{P} \left( \sup_{\zeta \in \mathcal{F}_\mathcal{W}, \zeta' \in \mathcal{F}_\theta} : |\langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_{\mu_N} - \langle \nabla_{\hat{Y}} \ell(g, f_\zeta), f_{\zeta'} \rangle_\mu | \geq \epsilon + B_{plr}(\mathcal{C}) \right) \quad (528)$$

$$\begin{aligned} &\leq \delta_{\mathcal{C}} + C \exp \left( \log \left( \mathcal{C}_{\mathcal{F}_W} \left( \frac{\epsilon}{8\tilde{L}_{\Phi} \max\{B_{\ell}, LB_{\Phi}\}} \right) \right) \right. \\ &\quad \left. + \log \left( \mathcal{C}_{\mathcal{F}_{\theta}} \left( \frac{\epsilon}{8\tilde{L}_{\Phi} \max\{B_{\ell}, LB_{\Phi}\}} \right) \right) - N \min \left\{ \frac{\epsilon^2}{4K^2}, \frac{\epsilon}{2K} \right\} \right). \end{aligned} \quad (529)$$

This completes our result.  $\square$

## E NUMERICAL EXPERIMENTS

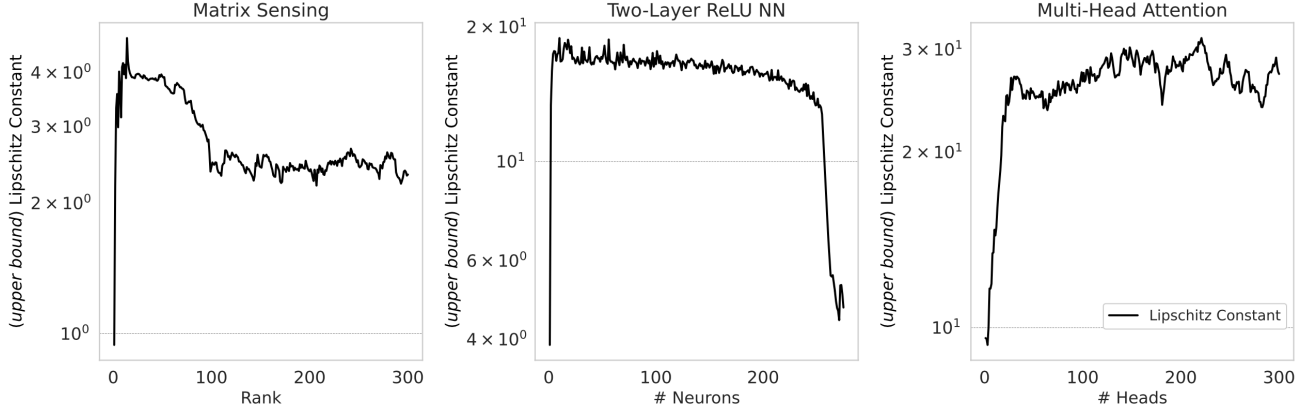


Figure 1: Numerical simulations of the Lipschitz constant (or upper bound thereof) obtained for different model widths ( $r$ ).

In this section, we present numerical simulations for the problems of low-rank matrix sensing, two-layer ReLU neural networks, and multi-head attention. In each simulation shown in Figure E, we generated data using a teacher model with random initialization of parameters,  $Y = \Phi_{r^*}(\{W_j\})(X) + \epsilon$ , where  $r^* = 64$ ,  $X \sim \mathcal{N}(0, (\sigma_X^2/n_X)I)$ , and  $\epsilon \in \mathcal{N}(0, (\sigma_E^2/n_Y)I)$ . We used gradient descent to reach a stationary point for each  $R$  (rank, number of neurons, or number of heads), starting from 1 and increasing up to 300. The first factor was initialized with small-scale random values. For each subsequent factor, we initialized the new factor with the supremum obtained from the polar equation (14), following the algorithm in Haeffele and Vidal (2015, 2020).

In each problem shown in Figure E, we plot the upper bounds on the Lipschitz constant for these problems. For matrix sensing, the Lipschitz constant is trivially upper-bounded by  $\|\mathbf{U}\mathbf{V}^T\|_2$ ; for the ReLU neural network, it is upper-bounded by  $\|\mathbf{U}\|_2\|\mathbf{V}\|_2$ ; and for multi-head attention, it is upper-bounded by  $\sum_{j=1}^r \|\mathbf{V}_j\|_2$ . We can observe from Figure E upper bounds on the Lipschitz constants are uniformly bounded, indicating that our Assumption 6 is realistic and holds empirically.

We conjecture that it is possible to show that the Lipschitz constants are uniformly bounded for any stationary. However, the analysis of this is beyond the scope of this work. Similar analyses based on gradient descent can be found in Oymak and Soltanolkotabi (2019).

## F OTHER RELATED WORKS

In this section, we provide a comprehensive study of the related works of the applications that are of the concern in this work.

**Statistical Learning Theory (SLT):** SLT provides a theoretical framework for analyzing generalization error, often producing results of the form (530). The seminal work by Vapnik (2000) established a systematic approach to deriving bounds of this nature. Over time, various approaches in SLT have attempted to estimate  $\epsilon(\mathcal{F}, N, \delta)$ , as summarized in Table 2. A recurring challenge in these bounds is the need to quantify the “capacity” of the model’s hypothesis class, which is particularly difficult for DNNs.

While there have been attempts to estimate the VC-dimension, such as those based on the norm of the parameters (Neyshabur et al., 2017), the resulting bounds heavily depend on the norm of the parameters. Consequently, it

remains unclear how to accurately estimate the sample complexity of models when varying the depth or width of DNNs. More recent work, such as Imaizumi and Schmidt-Hieber (2023), presents bounds that are tight but still dependent on the norm of the weights, assuming that the SGD iterates converge to a specific class of parameters.

Another line of research by Muthukumar and Sulam (2023) explores bounds that leverage the sparsity of feed-forward neural networks. However, there is still a lack of data-dependent bounds that do not rely on capacity estimates for models trained on random labels.

$$\mathbb{P} \left( \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{X,Y} [\ell(Y, f(X))] - \frac{1}{N} \sum_{i=1}^N \ell(Y_i, f(X_i)) \right| \leq \epsilon(\mathcal{F}, \delta, N) \right) \geq 1 - \delta \quad (530)$$

Description	$\epsilon(\mathcal{F}, N, \delta)$
Vapnik-Chernoviks Dimension, (Vapnik, 2000)	$\sqrt{\frac{\text{VCdim}(\mathcal{F}) - \log(\delta)}{N}}$
Rademacher Complexity, (Bartlett and Mendelson, 2001)	$R_N(\mathcal{F}) + \sqrt{\frac{-\log(\delta)}{N}}$
PAC-Bayes Bounds, (McAllester, 1999)	$\frac{KL(Q  P) - \log(\delta)}{N}$
Gaussian Complexity, (Bartlett and Mendelson, 2001)	$G_N(\mathcal{F}) + \sqrt{\frac{-\log(\delta)}{N}}$
Information-theoretic Bounds, (?)	$\frac{1}{N} \sum_{i=1}^N \sqrt{I(W; (X_i, Y_i))}$
Algorithmic Stability, (Feldman and Vondrak, 2019)	$\beta + \sqrt{\frac{-\log(\delta)}{N}}$

Table 2: SLT frameworks (in chronological order)

**Matrix recovery:** This is a fundamental problem in signal processing, where we seek to recover a matrix by indirect measurements, like random measurements, and random entry access. We typically have limited measurements; the problem itself is ill-posed when reconstructing the matrix. However, if the underlying matrix has certain special structures like low-rankedness, or sparsity in entries, the problem becomes tractable so as to reconstruct the true matrix. In practice, the problem tends to have low-rankedness, therefore having immense literature in this area, our work also presents such results, considering the optimization.

Let,  $Y_i = \langle M^*, X_i \rangle + \epsilon \in \mathbb{R}$ , where,  $X_i \in \mathbb{R}^{m \times n}$  ( $m \geq n$ ) is Gaussian entried matrix,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  and  $M^* \in \mathbb{R}^{m \times n}$  is a  $r^*$ -rank matrix. Consider the below problem

$\begin{aligned} & \min_{M \in \mathbb{R}^{m \times n}} \text{rank}(M) \\ \text{s.t.} \quad & \ Y_i - \langle M, X_i \rangle\  \leq \delta \end{aligned} \quad (531)$	$\begin{aligned} & \min_{M \in \mathbb{R}^{m \times n}} \ M\ _* \\ \text{s.t.} \quad & \ Y_i - \langle M, X_i \rangle\  \leq \delta \end{aligned} \quad (532)$
$\begin{aligned} & \min_{r \in \mathbb{N}, U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \ UV^T\ _* \\ \text{s.t.} \quad & \ Y_i - \langle UV^T, X_i \rangle\  \leq \delta \end{aligned} \quad (533)$	$\begin{aligned} & \min_{r \in \mathbb{N}, U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \frac{1}{2} (\ U\ _F^2 + \ V\ _F^2) \\ \text{s.t.} \quad & \ Y_i - \langle UV^T, X_i \rangle\  \leq \delta \end{aligned} \quad (534)$

Table 3: Optimization problems for matrix sensing

$$\min_{r \in \mathbb{N}, U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \|Y_i - \langle UV^T, X_i \rangle\|^2 + \frac{\lambda}{2} [\|U\|_F^2 + \|V\|_F^2] \quad (535)$$

The optimization problem in (531) is non-convex due to its rank-minimization nature, an NP-HARD problem. However, under certain specific conditions on the measurement matrices  $X_i$ , the convex relaxation (532) can recover solutions to (531), as demonstrated in Recht et al. (2008). Solving the convex program (532) requires computing the Singular Value Decomposition (SVD), which has a computational complexity of  $\mathcal{O}(mn^2)$ .

To mitigate this computational burden, the Burer-Monteiro (BM) factorization (Burer and Monteiro, 2003) is employed, yielding the bilinear factorization in the non-convex program (533). This approach is more efficient than (532) because it introduces an implicit rank constraint,  $\text{rank}(UV^T) \leq \min(n, r)$ , which reduces the runtime of SVD to  $\mathcal{O}((m+n)r^2)$ . Additionally, the equivalence between the nuclear norm and the sum of Frobenius norms,



as shown by Giampouras et al. (2020), further accelerates the optimization process, reducing the complexity to  $\mathcal{O}((m+n)r)$ .

While the BM factorization program (533) is non-convex, in contrast to the convex program (532), gradient descent (GD) algorithms typically guarantee only local minima for non-convex optimization problems (Reddy and Vidyasagar, 2023). However, Ge et al. (2017) has proven that the program (533) has no spurious local minima, and any local minimum is indeed a global minimum. Numerous studies (Jia et al., 2023) have explored the optimization landscapes and the convergence to global minima.

Our work primarily focuses on the generalization capabilities of the BM factorization program (535), which represents the Lagrangian form of the program (533). Table F summarizes the results from the literature that provide matrix recovery guarantees; from this we can suggest there are no bounds in the literature for low-rank matrix recovery with nuclear norm regularization under noisy settings with generic parameterization. Our work presents results first of its kind.

Measurement Type	Scenario	Reference	Result
<b>Exact</b>	Under-Parameterized ( $r < r^*$ )	N/A	N/A
	Exactly-Parameterized ( $r = r^*$ )	N/A	Not directly available.
	Over-Parameterized ( $r > r^*$ )	(Stöger and Soltanolkotabi, 2021)	$\ UU^T - M^*\ _F \lesssim r^{*1/8}(r - r^*)^{3/8}$ when $r \in (r^*, 2r^*)$ .
	Generic Parameterization ( $r \geq 1$ )	(Jin et al., 2023)	GD learns rank incrementally, $\ M^* - UU^T\ _F \lesssim \alpha^{c_2 k^2}$ , but analysis is algorithmic.
	SDP Relaxation ( <i>Full SDP Matrix</i> )	N/A	Not directly available.
<b>Noisy</b>	Under-Parameterized ( $r < r^*$ )	N/A	N/A
	Exactly-Parameterized ( $r = r^*$ )	(Ma et al., 2020)	$\ M^* - UU^T\ _F \lesssim \sqrt{\frac{\log(m)}{N}}$ under RIP assumptions $\delta_{4r^*} \leq 0.1$ .
		(Negahban and Wainwright, 2011)	$\ \hat{M} - M\ _F \lesssim \sqrt{r^* \frac{m+n}{N}}$ .
	Over-Parameterized ( $r > r^*$ )	(Ma et al., 2020)	$\ M^* - UU^T\ _F \lesssim \sqrt{\delta_{r+r^*}} \ M^*\ _2$ .
	Generic Parameterization ( $r \geq 1$ )	N/A	N/A
	SDP Relaxation ( <i>Full SDP Matrix</i> )	(Candès and Plan, 2011)	$\ \hat{M} - M^*\ _F \lesssim \sqrt{nr^*/N}$ under RIP assumptions.
		(Koltchinskii et al., 2011)	$\ \hat{M} - M\ _F \lesssim \frac{mnr^* \log(N)}{N}$ under uniform noisy measurements.

Table 4: Summary of Related Works on Matrix Recovery. N/A is an acronym for "Not Available".

**Transformers:** The remarkable success of Large Language Models (LLMs) (Team, 2024) can largely be attributed to their foundational architecture—Transformers (Vaswani et al., 2017). The optimization dynamics of Transformers have been a subject of extensive recent research (Bordelon et al., 2024), (Singh, 2023), (Yang et al., 2022), (Tian et al., 2023), (Nichani et al., 2024). Although Transformers exhibit impressive generalization capabilities in practical applications (Zhou et al., 2024), there is still a significant gap in the theoretical analysis of their generalization error.

To apply classical SLT bounds, one must determine the capacities of the function classes induced by Transformers. Previous attempts, such as in (Edelman et al., 2022), have made progress but were limited to scenarios where input data is bounded. In contrast, our work extends these results to settings where the inputs are not necessarily bounded.

Another line of research (Li et al., 2023), (Deora et al., 2024) has provided bounds that depend on step sizes and initialization choices for Gradient Descent (GD). For instance, Li et al. (2023) offered bounds within the context of in-context learning (Zhang et al., 2024), yet without evaluating the capacities of the stable algorithms used to train these Transformers.

In the broader literature, existing studies on generalization bounds often rely on strong assumptions, such as (i) bounded input data, (ii) algorithmic stability in some defined sense, and (iii) Lipschitz continuity of the loss function (which does not hold globally for mean squared error). Our results address these limitations by providing near-tight sample complexity bounds, offering a more comprehensive understanding of generalization in Transformer models.

## G PRELIMINARIES

This section provides preliminaries of convex analysis and concentration of measure.

### G.1 Convex Functions

**Definition 1** ( $L^2$  functions). *A function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is said to be square integrable on measure  $\mu$ , i.e.,  $L^2(\mu)$  if and only if,*

$$\langle f, f \rangle_\mu = \int_{x \in \mathcal{X}} \langle f(x), f(x) \rangle_{\mathcal{Y}} d\mu(x) < \infty. \quad (536)$$

**Definition 2** (Convex Set, (Rockafellar, 1970)). *A set  $\mathcal{C}$  is said to be convex if and only if  $\forall f, g \in \mathcal{C}, \alpha f + (1 - \alpha)g \in \mathcal{C}; \forall \alpha \in [0, 1]$ .*

**Definition 3** (Convex functions, (Rockafellar, 1970)). *A function,  $\Omega$  is said to be convex if and only if  $\text{dom}(\Omega)$  is convex and  $\forall f, g \in \text{dom}(\Omega)$  and any  $\alpha \in [0, 1]$ .*

$$\Omega(\alpha f + (1 - \alpha)g) \leq \alpha \Omega(f) + (1 - \alpha) \Omega(g). \quad (537)$$

**Definition 4** (Gauge function, (Rockafellar, 1970)). *The gauge function or the Minkowski functional is defined in a set  $\mathcal{C} \in L^2(\mu)$  for a point  $f$  as follows,*

$$\sigma_{\mathcal{C}}(f) := \inf \{t \geq 0; \text{ such that } f \in t \text{conv}(\mathcal{C})\}. \quad (538)$$

**Definition 5** (Polar Set, (Rockafellar, 1970)). *The polar set of any set  $\mathcal{C} \subseteq L^2(\mu)$  is given be*

$$\mathcal{C}^\circ := \{g \in L^2(\mu) : \text{ such that } \langle g, f \rangle_\mu \leq 1; \forall f \in \mathcal{C}\}. \quad (539)$$

**Proposition 5** (Polar Properties).

**Definition 6** (Polar function, (Rockafellar, 1970)). *The polar function of any gauge function,  $\sigma$  defined in the set  $\mathcal{C} \subseteq L^2(\mu)$  is given be*

$$\sigma_{\mathcal{C}}^\circ(g) := \sigma_{\mathcal{C}^\circ}(g). \quad (540)$$

**Definition 7** (Fenchel dual, (Rockafellar, 1970)). *The fenchel-dual for any  $\mu$ -measurable function,  $\Omega$  evaluated at  $g \in L^2(\mu)$  is defined by,*

$$\Omega^*(g) := \sup_{f \in L^2(\mu)} \langle g, f \rangle_\mu - \Omega(f). \quad (541)$$

**Lemma 12** (First Convexity, (Rockafellar, 1970)). *Any function  $\Omega$  that is first-order differentiable,  $\Omega \in \mathcal{C}^1$  is convex if and only if for any  $f, g \in \text{dom}(\Omega)$*

$$\Omega(f) \geq \Omega(g) + \langle \nabla \Omega(g), f - g \rangle_\mu. \quad (542)$$

**Lemma 13** (Strongly Convex, (Rockafellar, 1970)). *Any function  $\Omega$  that is first-order differentiable,  $\Omega \in \mathcal{C}^1$  is said to be  $\lambda(\geq 0)$ -strongly convex if and only if for any  $f, g \in \text{dom}(\Omega)$*

$$\Omega(f) \geq \Omega(g) + \langle \nabla \Omega(g), f - g \rangle_\mu + \frac{\lambda}{2} \|f - g\|_\mu^2. \quad (543)$$

**Definition 8** (Lipschitz Continuous). *A function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is said to be Lipschitz continuous with Lipschitz constant  $\|f\|_{\text{Lip}}$  if for any  $x_1, x_2 \in \mathcal{X}$*

$$\|f(x_1) - f(x_2)\|_{\mathcal{Y}} \leq \|f\|_{\text{Lip}} \|x_1 - x_2\|_{\mathcal{X}}. \quad (544)$$

**Remark:** Lipschitz constant,  $\|f\|_{\text{Lip}}$  is not a norm but only a semi-norm. Because  $\|f\|_{\text{Lip}} = 0$ , it implies that  $f$  can be any constant function.

**Definition 9** (Lipschitz Smooth). *A first-order differentiable function  $f : \mathcal{X} \rightarrow \mathcal{Y} \in \mathcal{C}^1$  is said to be Lipschitz smooth if  $\nabla f$  is Lipschitz continuous.*

**Definition 10** ( $(L, \lambda)$  convex function). *A first-order differentiable function  $f : \mathcal{X} \rightarrow \mathcal{Y} \in \mathcal{C}^1$  is said to be  $(L, \lambda)$  convex if and only if  $f$  is  $L$ -Lipschitz smooth and  $\lambda$ -strongly convex, here  $L \geq \lambda \geq 0$ .*

**Proposition 6** (Properties of Lipschitz). *The below are few properties of Lipschitz functions,*

1. *If function  $f : \mathcal{X} \rightarrow \mathcal{Y} \in \mathcal{C}^1$  then  $\sup_{x \in \mathcal{X}} \frac{\|\langle \nabla f(x), x \rangle\|_{\mathcal{Y}}}{\|x\|_{\mathcal{X}}} = \|f\|_{Lip}$ .*
2. *If convex function  $f : \mathcal{X} \rightarrow \mathcal{Y} \in \mathcal{C}^1$  is  $L$ -Lipschitz smooth then,*

$$f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle_{\mathcal{Y}} \leq f(x) \leq f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle_{\mathcal{Y}} + \frac{L}{2} \|x - x_0\|_{\mathcal{X}}^2. \quad (545)$$

3. *If convex function  $f : \mathcal{X} \rightarrow \mathcal{Y} \in \mathcal{C}^1$  is  $(L, \lambda)$  convex then,*

$$\begin{aligned} f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle_{\mathcal{Y}} + \frac{\lambda}{2} \|x - x_0\|_{\mathcal{X}}^2 &\leq f(x) \\ &\leq f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle_{\mathcal{Y}} + \frac{L}{2} \|x - x_0\|_{\mathcal{X}}^2. \end{aligned} \quad (546)$$

## G.2 Concentration of Measure

**Definition 11** (Greater than or approximately equal to). *The inequality  $f \gtrsim g$  means that  $\exists C > 0$  such that  $f \geq Cg$ .*

**Definition 12** (Sub-Gaussianity). *A random variable,  $X$  is said to be sub-Gaussian with proxy variance,  $\sigma^2$  if the following is satisfied,*

$$\mathbb{E}_X [e^t[X - \mathbb{E}[X]]] \leq \exp\left(-\frac{t^2\sigma^2}{2}\right); \forall t \geq 0. \quad (547)$$

We denote,  $X \sim SG(\sigma^2)$ .

**Definition 13** (Sub-exponential). *A random variable  $X$  is said to be subexponential with the proxy parameter  $\lambda$  if the following is satisfied*

$$\mathbb{E}_X [e^{t[X - \mathbb{E}[X]]}] \leq \exp\left(-\frac{t\lambda}{2}\right); \forall t \geq 0. \quad (548)$$

We denote  $X \sim SE(\lambda)$ .

**Proposition 7** (Properties of Sub-Gaussianity and Sub-exponential). *Let  $X, Y$  be two random variables that need not be independent.*

1.  *$X \in \mathbb{R}^n \sim SG\left(\frac{\sigma_X^2}{n} I_{n \times n}\right)$  if and only if  $\|X\|^2 \sim SE(\sigma_X^2)$ .*
2. *If  $X \in \mathbb{R}^n \sim SG\left(\frac{\sigma_X^2}{n} I_{n \times n}\right)$ , then for any Lipschitz function  $\phi : \mathcal{X} \rightarrow \mathbb{R}$ ,  $\phi(X) \sim SG(\|\phi\|_{Lip}^2 \sigma_X^2 / n)$ .*
3. *If  $X \in \mathbb{R}^n \sim SG\left(\frac{\sigma_X^2}{n} I_{n \times n}\right)$ , and  $Y \in \mathbb{R}^n \sim SG\left(\frac{\sigma_Y^2}{n} I_{n \times n}\right)$ , then  $\langle X, Y \rangle \sim SE(\sigma_X \sigma_Y)$ .*

**Lemma 14** (Uniform concentration of function). *Consider an  $n_X$ -dimensional vector  $X$ , and a parameterized function,  $g_\theta : \mathcal{X} \rightarrow \mathbb{R}$ , where  $\theta \in \mathcal{F}_\theta$ . Let  $\mathcal{C}$  be some convex set obeying  $P(\cap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}$ . Assume that for any fixed  $\theta_1, \theta_2 \in \mathcal{F}_\theta$  and any  $Z \in \mathcal{C}$  we have*

$$|g_{\theta_1}(Z) - g_{\theta_2}(Z)| \leq Kd(\theta_1, \theta_2). \quad (549)$$

*In addition, suppose that for any fixed  $\theta \in \mathcal{F}_\theta$ , we have*

$$|\mathbb{E}[g_\theta(\mathcal{P}_{\mathcal{C}}(X)) - g_\theta(X)]| \leq B, \quad (550)$$

where  $\mathcal{P}_C(\cdot)$  denotes the Euclidean projection onto the set  $\mathcal{C}$ . Finally, suppose that for any fixed  $\theta$  and  $\epsilon \in [-t, t]$  it holds that

$$P\left(\left|\int_{\omega} (g_{\theta} \circ \mathcal{P}_C) d\mu_N(\omega) - \int_{\omega} (g_{\theta} \circ \mathcal{P}_C) d\mu(\omega)\right| \geq \epsilon\right) \leq \delta(\epsilon), \quad (551)$$

Then for any  $\epsilon \in [-t, t]$ ,

$$P\left(\sup_{\theta \in \mathcal{F}_{\theta}} \left|\int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega)\right| \geq \epsilon + B\right) \leq \mathcal{N}(\mathcal{F}_{\theta}, d(\cdot, \cdot), \epsilon/(2K))\delta(\epsilon/4) + \delta_{\mathcal{C}}. \quad (552)$$

*Proof.* The proof technique is similar to that of (Li and Wei, 2023, Lemma 6) but includes more general parameter sets  $\mathcal{F}_{\theta}$ . Let us define

$$h_{\theta}(X) := g_{\theta}(\mathcal{P}_C(X)), \quad (553)$$

from the assumptions in the lemma, we have that,

$$P\left(\left|\int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega)\right| \geq \epsilon\right) \leq \delta(\epsilon), \quad (554)$$

Next, we must establish uniform concentration overall  $\theta \in \mathcal{E}_{\theta}$ . Let us construct a  $\nu$ -net for  $\mathcal{F}_{\theta}$ . For any  $\theta' \in \mathcal{N}_{\nu}(\mathcal{F}_{\theta}, d(\cdot, \cdot))$ ,  $\theta \in \mathcal{F}_{\theta}$  from the triangular inequality, and as

$$|h_{\theta} - h_{\theta'}| = |h_{\theta}(X) - h_{\theta'}(X)| = |g_{\theta}(\mathcal{P}_C(X)) - g_{\theta'}(\mathcal{P}_C(X))| \leq Kd(\theta, \theta'). \quad (555)$$

Then for any,  $X$  we have that that,

$$h_{\theta'} - Kd(\theta, \theta') \leq h_{\theta} \leq h_{\theta'} + Kd(\theta, \theta'). \quad (556)$$

Integrating with respect to the measure  $\mu_N$ , we obtain

$$\int_{\omega} h_{\theta'} d\mu_N(\omega) - Kd(\theta, \theta') \leq \int_{\omega} h_{\theta} d\mu_N(\omega) \leq \int_{\omega} h_{\theta'} d\mu_N(\omega) + L_X d(\theta, \theta'). \quad (557)$$

Similarly for the measure  $\mu$  we obtain

$$\implies \int_{\omega} h_{\theta'} d\mu(\omega) - Kd(\theta, \theta') \leq \int_{\omega} h_{\theta} d\mu(\omega) \leq \int_{\omega} h_{\theta'} d\mu(\omega) + Kd(\theta, \theta'). \quad (558)$$

Now, subtracting the above equations, we obtain

$$\begin{aligned} & \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) - 2Kd(\theta, \theta') \\ & \leq \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \leq \\ & \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) + 2Kd(\theta, \theta'). \end{aligned} \quad (559)$$

Now, take the absolute value on both sides. Later on, applying triangular inequality, we obtain

$$\left|\int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega)\right| \leq \left|\int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega)\right| + 2Kd(\theta, \theta'). \quad (560)$$

Now choose,  $\theta^*$  as  $\arg \sup_{\theta \in \mathcal{F}_{\theta}} \left|\int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega)\right|$ , then we have that,

$$\sup_{\theta \in \mathcal{F}_{\theta}} \left|\int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega)\right| \leq \left|\int_{\omega} h_{\theta^*} d\mu_N(\omega) - \int_{\omega} h_{\theta^*} d\mu(\omega)\right| + 2Kd(\theta^*, \theta'). \quad (561)$$

Now choose any  $\theta'$  that lies at-most  $\nu$  from  $\theta^*$  on the metric,  $d(.,.)$ , i.e,  $d(\theta', \theta^*) \leq \nu$ , we have

$$\sup_{\theta \in \mathcal{F}_\theta} \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| \leq \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right| + 2K\nu. \quad (562)$$

By definition, we can bound the right hand term by the supremum,

$$\sup_{\theta \in \mathcal{F}_\theta} \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| \leq 2K\nu + \sup_{\theta' \in \mathcal{N}_\nu(\mathcal{F}_\theta, d(.,.))} \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right|. \quad (563)$$

We apply the probability measure on both side, obtaining,

$$\begin{aligned} & P \left( \sup_{\theta \in \mathcal{F}_\theta} \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| \geq \epsilon \right) \\ & \leq P \left( \sup_{\theta' \in \mathcal{N}_\nu(\mathcal{F}_\theta, d(.,.))} \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right| \geq \epsilon - 2K\nu \right), \end{aligned} \quad (564)$$

the inequality is satisfied by the monotonicity of the probability measure. Now we apply the union-argument for the  $\nu$ -net cover then we have

$$\begin{aligned} & P \left( \sup_{\theta \in \mathcal{F}_\theta} \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| \geq \epsilon \right) \leq \\ & P \left( \bigcup_{\theta' \in \mathcal{N}_\nu(\mathcal{F}_\theta, d(.,.))} \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right| \geq \epsilon - 2K\nu \right), \end{aligned} \quad (565)$$

Now we upper bound the right side union term with summation, and then we have

$$\begin{aligned} & P \left( \sup_{\theta \in \mathcal{F}_\theta} \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| \geq \epsilon \right) \leq \\ & \sum_{\theta' \in \mathcal{N}_\nu(\mathcal{F}_\theta, d(.,.))} P \left( \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right| \geq \epsilon - 2K\nu \right), \end{aligned} \quad (566)$$

Now we replace the summation with the  $\nu$ -covering number,  $\mathcal{N}(\mathcal{F}_\theta, d(.,.), \nu)$  obtaining

$$P \left( \sup_{\theta \in \mathcal{F}_\theta} \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| \geq \epsilon \right) \leq \mathcal{N}(\mathcal{F}_\theta, d(.,.), \nu) \delta(\epsilon - 2K\nu). \quad (567)$$

Now set  $\nu = \epsilon/(2K)$  then we have

$$\implies P \left( \sup_{\theta \in \mathcal{F}_\theta} \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| \geq \epsilon \right) \leq \mathcal{N}(\mathcal{F}_\theta, d(.,.), \epsilon/(2K)) \delta(\epsilon/2). \quad (568)$$

Now, we have established the uniform concentration for  $h_\theta$ . Next, we move onto relating  $h_\theta$  with the desired function  $h_\theta$ .

Recall that

$$|\mathbb{E}[h_\theta(X) - g_\theta(X)]| \leq B. \quad (569)$$

As  $P(\cap_{i=1}^N X_i \in \mathcal{C}) \geq 1 - \delta_{\mathcal{C}}$ , we can safely claim  $\int_{\omega} g_\theta d\mu_N(\omega) = \int_{\omega} h_\theta d\mu_N(\omega)$  with probability at least  $1 - \delta_{\mathcal{C}}$ . We have

$$\left| \int_{\omega} g_\theta d\mu_N(\omega) - \int_{\omega} g_\theta d\mu(\omega) \right| = \left| \int_{\omega} h_\theta d\mu_N(\omega) - \int_{\omega} g_\theta d\mu(\omega) \right|$$

$$\leq \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| + \left| \int_{\omega} h_{\theta} d\mu(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right|, \quad (570)$$

$$\implies \left| \int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega) \right| \leq \left| \int_{\omega} h_{\theta} d\mu_N(\omega) - \int_{\omega} h_{\theta} d\mu(\omega) \right| + B, \quad (571)$$

with probability at least  $1 - \delta_{\mathcal{C}}$ . Now we check the Lipschitzness of function  $h_{\theta}$  in  $\theta$  we have

$$|h_{\theta_1}(X) - h_{\theta_2}(X)| = |g_{\theta_1}(\mathcal{P}_{\mathcal{C}}(X)) - g_{\theta_2}(\mathcal{P}_{\mathcal{C}}(X))| \leq Kd(\theta_1, \theta_2). \quad (572)$$

Similarly, in expectation measure, we have that

$$\begin{aligned} |\mathbb{E}[h_{\theta_1}(X)] - \mathbb{E}[h_{\theta_2}(X)]| &= |\mathbb{E}[g_{\theta_1}(\mathcal{P}_{\mathcal{C}}(X))] - \mathbb{E}[g_{\theta_2}(\mathcal{P}_{\mathcal{C}}(X))]| \leq \mathbb{E}[|g_{\theta_1}(\mathcal{P}_{\mathcal{C}}(X)) - g_{\theta_2}(\mathcal{P}_{\mathcal{C}}(X))|] \\ &\leq Kd(\theta_1, \theta_2). \end{aligned} \quad (573)$$

Consequently for any  $\theta' \in \{\theta' : d(\theta, \theta') \leq \epsilon/(2K)\}$ , we have that

$$\left| \int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega) \right| \leq \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right| + \nu + B. \quad (574)$$

Now we choose  $\theta = \theta^* = \sup_{\theta \in \mathcal{F}_{\theta}} \left| \int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega) \right|$  then,

$$\begin{aligned} \sup_{\theta \in \mathcal{F}_{\theta}} \left| \int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega) \right| &= \left| \int_{\omega} g_{\theta^*} d\mu_N(\omega) - \int_{\omega} g_{\theta^*} d\mu(\omega) \right| \\ &\leq \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right| + \epsilon + B. \end{aligned} \quad (575)$$

We can take a supremum over  $\theta'$  in the upper bound of the right side term; we have

$$\sup_{\theta \in \mathcal{F}_{\theta}} \left| \int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega) \right| \leq B + \epsilon + \sup_{\theta' \in \mathcal{F}_{\theta'}} \left| \int_{\omega} h_{\theta'} d\mu_N(\omega) - \int_{\omega} h_{\theta'} d\mu(\omega) \right|. \quad (576)$$

Then we use the inequality (568) and (571) we have that,

$$\sup_{\theta \in \mathcal{F}_{\theta}} \left| \int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega) \right| \leq 2\epsilon + B, \quad (577)$$

with probability at least  $1 - [\mathcal{N}(\mathcal{F}_{\theta}, d(\cdot, \cdot), \epsilon/(2K))\delta(\epsilon/2) + \delta_{\mathcal{C}}]$ . Now rescaling we obtain that

$$P \left( \sup_{\theta \in \mathcal{F}_{\theta}} \left| \int_{\omega} g_{\theta} d\mu_N(\omega) - \int_{\omega} g_{\theta} d\mu(\omega) \right| \geq \epsilon + B \right) \leq \mathcal{N}(\mathcal{F}_{\theta}, d(\cdot, \cdot), \epsilon/(2K))\delta(\epsilon/4) + \delta_{\mathcal{C}}. \quad (578)$$

□