

Learning Multi-Index Models with Neural Networks via Mean-Field Langevin Dynamics

Alireza Mousavi-Hosseini

University of Toronto and Vector Institute

MOUSAVI@CS.TORONTO.EDU

Denny Wu

New York University and Flatiron Institute

DENNYWU@NYU.EDU

Murat A. Erdogdu

University of Toronto and Vector Institute

ERDOGDU@CS.TORONTO.EDU

Abstract

We study the problem of learning multi-index models in high-dimensions using a two-layer neural network trained with the mean-field Langevin algorithm. Under mild distributional assumptions on the data, we characterize the *effective dimension* d_{eff} that controls both sample and computational complexity by utilizing the adaptivity of neural networks to latent low-dimensional structures. When the data exhibit such a structure, d_{eff} can be significantly smaller than the ambient dimension. We prove that the sample complexity grows linearly with d_{eff} , bypassing the limitations of the information exponent or the leap complexity that appeared in recent analyses of gradient-based feature learning. On the other hand, the computational complexity may inevitably grow exponentially with d_{eff} in the worst-case scenario.

1. Introduction

A key characteristic of neural networks is their adaptability to the underlying statistical model. Several works have shown that shallow neural networks trained by (variants of) gradient descent can efficiently learn functions of low-dimensional projections (i.e., multi-index models) with a sample complexity that depends on the properties of the nonlinear link function known as the *information exponent* [9] or the *leap complexity* [2]. Specifically, to learn a target function with information exponent or leap complexity $k \in \mathbb{N}_+$ on isotropic Gaussian data, a sample size of $n \gtrsim d^{\Theta(k)}$ is typically required in these analyses [1, 11, 12, 21, 35]. This sample complexity is also predicted by the framework of correlational statistical query (CSQ) lower bounds [2, 23].

On the other hand, if the (polynomial) optimization budget is not taken into consideration, [7] showed that neural networks can learn multi-index models with a sample complexity that does not depend on the information or leap exponent. However, thus far it has been relatively unclear whether standard first-order optimization algorithms inherit this optimality.

A promising approach to obtain statistically optimal sample complexity is to consider training neural networks in the *mean-field regime* [17, 34, 37, 41, 43], where overparameterization lifts the gradient descent dynamics into the space of measures with global convergence guarantees. While most existing results in this regime focus on optimization instead of generalization/learnability guarantees, recent works have shown that in certain settings, neural networks in the mean-field regime can achieve a sample complexity that does not depend on the leap complexity [18, 45, 47, 51]. However, these guarantees rely on stringent assumptions on data (isotropic Gaussian, hypercube, etc.) as well as single-index models with specific link functions [10, 33], or k -sparse party classifi-

cation [45, 47, 51]. Moreover, the computational complexity of the training algorithm in [45, 47] is exponential in the ambient (input) dimension, without adapting to the potential structure in the inputs. An exception in this direction is the recent work of [39], which considers the k -parity problem on anisotropic data.

Our contributions. Motivated by the above discussion, we address two key questions. First, *Can we train two-layer neural networks using the MFLA to learn arbitrary multi-index models with an optimal sample complexity?*

We answer this in the affirmative by showing that empirical risk minimization on a standard variant of a two-layer neural network can be achieved by the MFLA. This result handles arbitrary multi-index models on subGaussian data with general covariance, hence enabling us to achieve the optimal sample complexity with standard gradient-based training. However, such a universal guarantee will inevitably suffer from an exponential computational complexity; thus, the second fundamental question we aim to answer is

Are there conditions under which the computational complexity of the MFLA can be improved by adapting to data structure?

We provide a positive answer by showing that the computational complexity of MFLA is governed by the *effective dimension* of the learning problem, instead of its ambient dimension; this implies an improved efficiency of MFLA when the data is anisotropic as in prior works [28, 36].

Related works. The training dynamics of neural networks in the mean-field regime is described by a nonlinear partial differential equation in the space of parameter distributions [17, 34, 41]. Unlike the NTK description [19, 30] that freezes the parameters around the random initialization, the mean-field regime allows for the parameters to travel and learn useful features, leading to improved statistical efficiency. While convergence analyses for mean-field neural networks are typically *qualitative* in nature, in that they do not specify the rate of convergence or finite-width discrepancy, the mean-field Langevin algorithm that we study is a noticeable exception, for which the convergence rate [16, 29, 38] as well as uniform-in-time propagation of chaos [13, 44] have been established.

The benefit of feature learning has also been studied in a “narrow-width” setting for learning low-dimensional target functions. Examples of low-dimensional targets include single-index models [5, 9, 11, 21, 35] and multi-index models [1, 2, 12, 20, 23, 24]. While the information-theoretic threshold for learning such functions is $n \gtrsim d$ [8, 22], the complexity of gradient-based learning is governed by properties of the link function. For instance, in the single-index setting, prior works established a sufficient sample size of $n \gtrsim d^{\Theta(k)}$ where k is the *information exponent* [9, 11, 21] or the *generative exponent* [4, 22, 25, 32]. This presents a gap between the information-theoretically achievable sample complexity and the performance of neural networks optimized by gradient descent, which we aim to close by studying the statistical efficiency of mean-field neural networks.

2. Preliminaries: Statistical Model and Training Algorithm

Statistical model. In this paper, we consider the regression setting where the input $\mathbf{x} \in \mathbb{R}^d$ is generated from some distribution and the response $y \in \mathbb{R}$ is given by the multi-index model

$$y = g\left(\frac{\langle \mathbf{u}_1, \mathbf{x} \rangle}{\sqrt{k}}, \dots, \frac{\langle \mathbf{u}_k, \mathbf{x} \rangle}{\sqrt{k}}\right) + \xi. \quad (2.1)$$

Here, $g : \mathbb{R}^k \rightarrow \mathbb{R}$ is the unknown link function, ξ is a zero-mean ς -subGaussian noise independent from \mathbf{x} ; for simplicity, we assume that $\varsigma^2 \lesssim 1$. Without loss of generality, we assume $\mathbf{u}_1, \dots, \mathbf{u}_k$

are orthonormal and form the matrix $\mathbf{U} = (\mathbf{u}_1/\sqrt{k}, \dots, \mathbf{u}_k/\sqrt{k})^\top \in \mathbb{R}^{k \times d}$; thus, we can use the shorthand notation $y = g(\mathbf{U}\mathbf{x}) + \xi$. Throughout the paper, we consider the setting $k \ll d$, and treat k as an absolute constant independent from the ambient input dimension d .

For our predictor, we use a two-layer neural network coupled with ℓ_2 regularization to learn the statistical model (2.1). Denoting the m neurons with a matrix $\mathbf{W} := (\mathbf{w}_1, \dots, \mathbf{w}_m)^\top$, the student model and the ℓ_2 -regularizer are given as

$$\hat{y}_m(\mathbf{x}; \mathbf{W}) := \frac{1}{m} \sum_{j=1}^m \Psi(\mathbf{x}; \mathbf{w}_j) \quad \text{and} \quad R(\mathbf{W}) := \frac{1}{m} \|\mathbf{W}\|_{\text{F}}^2 = \frac{1}{m} \sum_{j=1}^m \|\mathbf{w}_j\|^2, \quad (2.2)$$

where $\Psi : \mathbb{R}^d \times \mathcal{W} \rightarrow \mathbb{R}$ is the activation function, and $\mathbf{w}_j \in \mathcal{W}$ with \mathcal{W} denoting the weight space. In this formulation, the second layer weights are all fixed to be +1. In the classical regression setting where we observe n i.i.d. samples $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ from the data distribution, the regularized population and empirical risks are defined respectively as

$$J_\lambda(\mathbf{W}) := \mathbb{E}[\ell(\hat{y}_m(\mathbf{x}; \mathbf{W}), y)] + \frac{\lambda}{2} R(\mathbf{W}) \quad \text{and} \quad \hat{J}_\lambda(\mathbf{W}) := \frac{1}{n} \sum_{i=1}^n \ell(\hat{y}_m(\mathbf{x}^{(i)}; \mathbf{W}), y^{(i)}) + \frac{\lambda}{2} R(\mathbf{W}),$$

where $\ell(\hat{y}, y) = \rho(\hat{y} - y)$ with $\rho : \mathbb{R} \rightarrow \mathbb{R}_+$ being a convex loss.

Training Algorithm. We minimize the regularized empirical risk $\hat{J}_\lambda(\mathbf{W})$ via the *mean-field Langevin algorithm* (MFLA) with stepsize η , which updates the weights at iteration l by

$$\mathbf{w}_j^{l+1} = \mathbf{w}_j^l - m\eta \nabla_{\mathbf{w}_j} \hat{J}_\lambda(\mathbf{W}) + \sqrt{\frac{2\eta}{\beta}} \boldsymbol{\xi}_j^l, \quad 1 \leq j \leq m, \quad (2.3)$$

where $\boldsymbol{\xi}_j^l$ are independent standard Gaussian random vectors. In Appendix A, we will introduce the necessary background on the mean-field Langevin dynamics (MFLD) and optimization on the space of measures, where we observe that (2.3) is a simple time-discretization of the MFLD.

3. Learning Multi-index Models

In this section, we will provide the learning guarantees. For technical reasons, we use an approximation of ReLU denoted by $z \mapsto \phi_{\kappa, \iota}(z)$ for some $\kappa, \iota > 1$, which is given by $\phi_{\kappa, \iota}(z) = \kappa^{-1} \ln(1 + \exp(\kappa z))$ for $z \in (-\infty, \iota/2]$ and extended on (ι, ∞) such that $\phi_{\kappa, \iota}$ is C^2 smooth, $|\phi_{\kappa, \iota}| \leq \iota$, $|\phi'_{\kappa, \iota}| \leq 1$, and $|\phi''_{\kappa, \iota}| \leq \kappa$. Note that $\phi_{\kappa, \iota}$ recovers ReLU as $\kappa, \iota \rightarrow \infty$. To be able to learn functions with positive and negative parts, we choose $\mathcal{W} = \mathbb{R}^{2d+2}$, and use the notation $\mathbf{w} = (\boldsymbol{\omega}_1^\top, \boldsymbol{\omega}_2^\top)^\top$ with $\boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \in \mathbb{R}^{d+1}$, and ultimately use

$$\Psi(\mathbf{x}; \mathbf{w}) := \phi_\kappa(\langle \tilde{\mathbf{x}}, \boldsymbol{\omega}_1 \rangle) - \phi_\kappa(\langle \tilde{\mathbf{x}}, \boldsymbol{\omega}_2 \rangle), \quad (3.1)$$

where $\tilde{\mathbf{x}} := (\mathbf{x}, \tilde{r}_x)^\top \in \mathbb{R}^{d+1}$ for a constant \tilde{r}_x corresponding to bias, to be specified later. The above can also be viewed as a 2-layer neural network with second-layer weights frozen at ± 1 . We make the following assumption on the input distribution.

Assumption 1 *The input \mathbf{x} has zero mean and covariance $\boldsymbol{\Sigma}$. Further, $\|\mathbf{x}\|$ and $\|\mathbf{U}\mathbf{x}\|$ are sub-Gaussian with respective norms $\sigma_n \|\boldsymbol{\Sigma}\|_{\text{F}}^{1/2}$ and $\sigma_u \|\boldsymbol{\Sigma}^{1/2} \mathbf{U}^\top\|_{\text{F}}$ for some absolute constants σ_n, σ_u .*

Even though the above assumption covers a wide range of input distributions, it is mainly motivated by the Gaussian case which satisfies the above assumption. Without loss of generality, we will consider a scaling where $\|\Sigma\| \lesssim 1$. A key quantity in our analysis is the *effective dimension* which governs the algorithmic guarantees.

Definition 1 (Effective dimension) Define $d_{\text{eff}} := c_x^2/r_x^2$ where $c_x := \text{tr}(\Sigma)^{1/2}$, $r_x := \|\Sigma^{1/2}U^\top\|_{\text{F}}$.

The effective dimension d_{eff} can be significantly smaller than the ambient dimension d , leading to particularly favorable results in the following when $d_{\text{eff}} = \text{polylog}(d)$. This concept has numerous applications from learning theory to statistical estimation; see e.g. [6, 27, 48, 50]. We make the following assumption on the link function in (2.1).

Assumption 2 g is locally Lipschitz, i.e. $|g(\mathbf{z}_1) - g(\mathbf{z}_2)| \leq L\|\mathbf{z}_1 - \mathbf{z}_2\|$ for $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^k$ satisfying $\|\mathbf{z}_1\| \vee \|\mathbf{z}_2\| \leq \tilde{r}_x := r_x(1 + \sigma_u\sqrt{2\ln(n)^q})$ for some $q > 0$ and $L = \mathcal{O}(1/r_x)$.

We emphasize that the above Lipschitz condition is only local; allowing e.g. polynomially growing link functions. We scale the Lipschitz constant with $1/r_x$ to make sure $|y| \asymp 1$ with high probability. The main result of this section is stated in the following theorem.

Theorem 2 For an appropriate choice of hyperparameters $\eta, \kappa, \iota, \lambda$, and β , with a sufficient number of samples, number of neurons, and number of iterations that can be bounded by

$$n \leq \tilde{\mathcal{O}}(d_{\text{eff}}), \quad m \leq \tilde{\mathcal{O}}(d^2 e^{\tilde{\mathcal{O}}(d_{\text{eff}})}), \quad l \leq \tilde{\mathcal{O}}(d^3 e^{\tilde{\mathcal{O}}(d_{\text{eff}})}), \quad (3.2)$$

with probability at least $1 - \mathcal{O}(n^{-q})$ for some $q > 0$, MFLA can achieve the excess risk bound

$$\mathbb{E}_{\mathbf{W}^l} \mathbb{E}_{y, \mathbf{x}}[\rho(y - \hat{y}_m(\mathbf{x}; \mathbf{W}^l))] - \mathbb{E}_\xi[\rho(\xi)] \leq o_n(1). \quad (3.3)$$

We refer to Theorem 28 in Appendix B for a more precise statement with the choice of hyperparameters. The theorem above demonstrates a certain adaptivity to the effective low-dimensional structure, both in terms of *statistical* and *computational* complexity, which occurs without explicitly encoding any information about the covariance structure in the algorithm. In contrast, “fixed-grid” methods (see [15] and references therein), that fix the first-layer of a two-layer network’s representation similar to random features regression [40], and then train the second-layer by solving a convex problem, do not show this type of adaptivity to low dimensions. In particular their computational complexity always scales exponentially with the ambient dimension d , unless information about the covariance structure is explicitly used when specifying the fixed representation.

Comparison with prior bounds. Here, we compare the guarantee of Theorem 2 with two prior works that are particularly relevant. First, [7] requires $d^{\frac{k+3}{2}}$ sample complexity for learning general multi-index models with k indices, which is worse than the complexity d_{eff} of Theorem 2 even in the worst case $d_{\text{eff}} = d$. The improvement in our bound is due to a refined control over $\|\mathbf{U}\mathbf{x}\|$. Further, [7] does not provide a quantitative analysis of the optimization complexity, and it is not clear if their algorithm is adaptive to the covariance structure. Moreover, [39] studies learning k -sparse parities, a subclass of multi-index models we considered, for which it is considerably simpler to construct optimal neural networks with a bounded activation. While the effective dimension (and the resulting sample complexity) of [39] is not explicitly scale-invariant, we derive a scale-invariant translation of their bound in Appendix C, and show that it is always lower bounded by our effective dimension, especially when Σ is nearly rank-deficient.

4. Interpreting the Effective Dimension

To better demonstrate the impact of effective dimension, we consider two covariance models.

Spiked covariance. We consider the spiked covariance model of [36]. Namely, given a spike direction $\boldsymbol{\theta} \in \mathbb{S}^{d-1}$, suppose the covariance and the target directions satisfy

$$\boldsymbol{\Sigma} = \frac{\mathbf{I}_d + \alpha \boldsymbol{\theta} \boldsymbol{\theta}^\top}{1 + \alpha}, \quad \alpha \asymp d^{\gamma_2}, \quad \|\mathbf{U} \boldsymbol{\theta}\| \asymp d^{-\gamma_1}, \quad \gamma_2 \in [0, 1], \quad \gamma_1 \in [0, 1/2]. \quad (4.1)$$

Note that in high-dimensional settings, $\gamma_1 = 1/2$ corresponds to a regime where $\boldsymbol{\theta}$ is sampled uniformly over \mathbb{S}^{d-1} , whereas $\gamma_1 = 0$ corresponds to the case where $\boldsymbol{\theta}$ has a strong correlation with \mathbf{U} . We only consider $\gamma_2 \leq 1$ since $\gamma_2 > 1$ corresponds to a setting where the input is effectively one-dimensional. In this setting, effective dimension depends on γ_1 and γ_2 .

Corollary 3 *Under the spiked covariance model (4.1), we have $d_{\text{eff}} \asymp d^{1 - \{(\gamma_2 - 2\gamma_1) \vee 0\}}$.*

To get improvements over the isotropic effective dimension d , either the spike magnitude α or the spike-target alignment $\|\mathbf{U} \boldsymbol{\theta}\|$ needs to be sufficiently large so that $\gamma_2 > 2\gamma_1$. As $\gamma_2 \rightarrow 1$ and $\gamma_1 \rightarrow 0$, the effective dimension will be smaller than $\text{polylog}(d)$, leading to a computational complexity that is quasipolynomial in d .

Covariance with decaying eigenvalues. Next, we consider a more general power-law decay for the eigenspectrum. Suppose $\boldsymbol{\Sigma} = \sum_{i=1}^d \lambda_i \boldsymbol{\theta}_i \boldsymbol{\theta}_i^\top$ is the spectral decomposition of $\boldsymbol{\Sigma}$, and

$$\frac{\lambda_i}{\lambda_1} \asymp i^{-\alpha}, \quad \frac{\|\mathbf{U} \boldsymbol{\theta}_i\|^2}{\|\mathbf{U} \boldsymbol{\theta}_1\|^2} \asymp i^{-\gamma}, \quad \text{for } 1 \leq i \leq d, \quad (4.2)$$

for some absolute constants $\alpha, \gamma > 0$. Notice that $\sum_{i=1}^d \|\mathbf{U} \boldsymbol{\theta}_i\|^2 = \|\mathbf{U}\|_{\text{F}}^2 = 1$. The following corollary characterizes d_{eff} in terms of the parameters α and γ .

Corollary 4 *Under the power-law eigenspectrum for the covariance matrix (4.2), we have*

$$d_{\text{eff}} \asymp \begin{cases} d^{1 \wedge (2 - \alpha - \gamma)} & \alpha < 1, \gamma < 1 \\ d^{1 - \alpha} & \alpha < 1, \gamma \geq 1, \\ d^{(1 - \gamma) \vee 0} & \alpha \geq 1 \end{cases}$$

where \asymp above hides $\text{polylog}(d)$ dependencies.

Thus, the computational complexity becomes quasipolynomial in d when $\alpha, \gamma \geq 1$. This happens when $\boldsymbol{\Sigma}$ is approximately low-rank with most of its eigenspectrum concentrated around the first few eigenvalues, and the corresponding eigenvectors are well-aligned with the row space of \mathbf{U} .

5. Conclusion

In this paper, we investigated the mean-field Langevin algorithm for learning multi-index models. We proved that the statistical and computational complexity of this problem can be characterized by an effective dimension which captures the low-dimensional structure in the input covariance, along with its correlation with the target directions. In particular, the sample complexity scales almost linearly with the effective dimension, while without additional assumptions, the computational complexity may scale exponentially with this quantity. We leave open the question of under which assumptions the computational complexity will be polynomial even in the effective dimension as an interesting direction for future research.

References

- [1] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. The merged-staircase property: a necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks. In *Conference on Learning Theory*, 2022.
- [2] Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz. Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics. *arXiv preprint arXiv:2302.11055*, 2023.
- [3] Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2005.
- [4] Luca Arnaboldi, Yatin Dandi, Florent Krzakala, Luca Pesce, and Ludovic Stephan. Repetita iuvant: Data repetition allows sgd to learn high-dimensional multi-index functions. *arXiv preprint arXiv:2405.15459*, 2024.
- [5] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation. *arXiv preprint arXiv:2205.01445*, 2022.
- [6] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the presence of low-dimensional structure: a spiked random matrix perspective. *Advances in Neural Information Processing Systems*, 36, 2023.
- [7] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal of Machine Learning Research*, 18(1):629–681, 2017.
- [8] Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Optimal errors and phase transitions in high-dimensional generalized linear models. *Proceedings of the National Academy of Sciences*, 116(12):5451–5460, 2019.
- [9] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *J. Mach. Learn. Res.*, 22:106–1, 2021.
- [10] Raphaël Berthier, Andrea Montanari, and Kangjie Zhou. Learning time-scales in two-layers neural networks. *arXiv preprint arXiv:2303.00055*, 2023.
- [11] Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with shallow neural networks. In *Advances in Neural Information Processing Systems*, 2022.
- [12] Alberto Bietti, Joan Bruna, and Loucas Pillaud-Vivien. On learning gaussian multi-index models with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023.
- [13] Fan Chen, Zhenjie Ren, and Songbo Wang. Uniform-in-time propagation of chaos for mean field langevin dynamics. *arXiv preprint arXiv:2212.03050*, 2022.
- [14] Zixiang Chen, Yuan Cao, Quanquan Gu, and Tong Zhang. A generalized neural tangent kernel analysis for two-layer neural networks. *Advances in Neural Information Processing Systems*, 33:13363–13373, 2020.

- [15] Lénaïc Chizat. Convergence rates of gradient methods for convex optimization in the space of measures. *Open Journal of Mathematical Optimization*, 3:1–19, 2022.
- [16] Lénaïc Chizat. Mean-field langevin dynamics: Exponential convergence and annealing. *arXiv preprint arXiv:2202.01009*, 2022.
- [17] Lenaic Chizat and Francis Bach. On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport. In *Advances in Neural Information Processing Systems*, 2018.
- [18] Lénaïc Chizat and Francis Bach. Implicit Bias of Gradient Descent for Wide Two-layer Neural Networks Trained with the Logistic Loss. In *Conference on Learning Theory*, 2020.
- [19] Lenaic Chizat, Edouard Oyallon, and Francis Bach. On Lazy Training in Differentiable Programming. In *Advances in Neural Information Processing Systems*, 2019.
- [20] Elizabeth Collins-Woodfin, Courtney Paquette, Elliot Paquette, and Inbar Seroussi. Hitting the high-dimensional notes: An ode for sgd learning dynamics on glms and multi-index models. *arXiv preprint arXiv:2308.08977*, 2023.
- [21] Alex Damian, Eshaan Nichani, Rong Ge, and Jason D Lee. Smoothing the landscape boosts the signal for sgd: Optimal sample complexity for learning single index models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [22] Alex Damian, Loucas Pillaud-Vivien, Jason D Lee, and Joan Bruna. The computational complexity of learning gaussian single-index models. *arXiv preprint arXiv:2403.05529*, 2024.
- [23] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural Networks can Learn Representations with Gradient Descent. In *Conference on Learning Theory*, 2022.
- [24] Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. Learning two-layer neural networks, one (giant) step at a time. *arXiv preprint arXiv:2305.18270*, 2023.
- [25] Yatin Dandi, Emanuele Troiani, Luca Arnaboldi, Luca Pesce, Lenka Zdeborová, and Florent Krzakala. The benefits of reusing batches for gradient descent in two-layer networks: Breaking the curse of information and leap exponents. *arXiv preprint arXiv:2402.03220*, 2024.
- [26] Monroe D Donsker and SR Srinivasa Varadhan. Asymptotic evaluation of certain markov process expectations for large time. iv. *Communications on pure and applied mathematics*, 36 (2):183–212, 1983.
- [27] B. Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of Lazy Training of Two-layers Neural Networks. In *Advances in Neural Information Processing Systems*, 2019.
- [28] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When Do Neural Networks Outperform Kernel Methods? In *Advances in Neural Information Processing Systems*, 2020.
- [29] Kaitong Hu, Zhenjie Ren, David Siska, and Lukasz Szpruch. Mean-field langevin dynamics and energy landscape of neural networks. *arXiv preprint arXiv:1905.07769*, 2019.

- [30] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. In *Advances in Neural Information Processing Systems*, 2018.
- [31] Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker-planck equation. *SIAM journal on mathematical analysis*, 29(1):1–17, 1998.
- [32] Jason D. Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional polynomials with sgd near the information-theoretic limit. *arXiv preprint arXiv:2406.01581*, 2024.
- [33] Arvind Mahankali, Haochen Zhang, Kefan Dong, Margalit Glasgow, and Tengyu Ma. Beyond ntk with vanilla gradient descent: A mean-field analysis of neural networks with polynomial width, samples, and time. *Advances in Neural Information Processing Systems*, 36, 2023.
- [34] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33): E7665–E7671, 2018.
- [35] Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A Erdogdu. Neural networks efficiently learn low-dimensional representations with SGD. In *The Eleventh International Conference on Learning Representations*, 2023.
- [36] Alireza Mousavi-Hosseini, Denny Wu, Taiji Suzuki, and Murat A Erdogdu. Gradient-based feature learning under structured data. *Advances in Neural Information Processing Systems*, 36, 2023.
- [37] Atsushi Nitanda and Taiji Suzuki. Stochastic particle gradient descent for infinite ensembles. *arXiv preprint arXiv:1712.05438*, 2017.
- [38] Atsushi Nitanda, Denny Wu, and Taiji Suzuki. Convex analysis of the mean field langevin dynamics. In *International Conference on Artificial Intelligence and Statistics*, pages 9741–9757. PMLR, 2022.
- [39] Atsushi Nitanda, Kazusato Oko, Taiji Suzuki, and Denny Wu. Improved statistical and computational complexity of the mean-field langevin dynamics under structured data. In *The Twelfth International Conference on Learning Representations*, 2024.
- [40] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, 2007.
- [41] Grant M Rotskoff and Eric Vanden-Eijnden. Neural networks as Interacting Particle Systems: Asymptotic convexity of the Loss Landscape and Universal Scaling of the Approximation Error. *arXiv preprint arXiv:1805.00915*, 2018.
- [42] Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94, 2015.
- [43] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A law of large numbers. *SIAM Journal on Applied Mathematics*, 80(2):725–752, 2020.

- [44] Taiji Suzuki, Denny Wu, and Atsushi Nitanda. Convergence of mean-field langevin dynamics: Time and space discretization, stochastic gradient, and variance reduction. *arXiv preprint arXiv:2306.07221*, 2023.
- [45] Taiji Suzuki, Denny Wu, Kazusato Oko, and Atsushi Nitanda. Feature learning via mean-field langevin dynamics: classifying sparse parities and beyond. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [46] Alain-Sol Sznitman. Topics in propagation of chaos. *Lecture notes in mathematics*, pages 165–251, 1991.
- [47] Matus Telgarsky. Feature selection and low test error in shallow low-rotation relu networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- [48] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.
- [49] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- [50] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.
- [51] Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and optimization of neural nets vs their induced kernel. *Advances in Neural Information Processing Systems*, 32, 2019.

Appendix A. Background on Optimization in Measure Space and MFLD

Notation. We denote the Euclidean inner product with $\langle \cdot, \cdot \rangle$, the Euclidean norm for vectors and the operator norm for matrices with $\|\cdot\|$, and the Frobenius norm with $\|\cdot\|_F$. We use $\mathcal{P}(\mathcal{W})$, $\mathcal{P}_2(\mathcal{W})$, and $\mathcal{P}_2^{ac}(\mathcal{W})$ to denote the set of (Borel) probability measures, the set of probability measures with finite second moment, and the set of absolutely continuous probability measures with finite second moment on the weight space \mathcal{W} , respectively. Finally, $\delta_{\mathbf{w}_0}$ denotes the Dirac measure at \mathbf{w}_0 .

Optimization in measure space. To minimize the regularized empirical risk \hat{J}_λ defined in Section 2, we will consider a discretization of the following set of SDEs, which essentially define an interacting particle system over m neurons:

$$d\mathbf{w}_j^t = -m\nabla_{\mathbf{w}_j} J(\mathbf{w}_1^t, \dots, \mathbf{w}_m^t)dt + \sqrt{\frac{2}{\beta}}d\mathbf{B}_t^j \quad \text{for } 1 \leq j \leq m, \quad (\text{A.1})$$

where $(\mathbf{B}_t^j)_{j=1}^m$ is a set of independent Brownian motions on the weight space \mathcal{W} .

Notice that the neural network and the regularizer in (2.2) are both invariant under permutations of the weights $(\mathbf{w}_1, \dots, \mathbf{w}_m)$; thus, an equivalent integral representation is given by

$$\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) := \int \Psi(\mathbf{x}; \cdot) d\mu_{\mathbf{W}} \quad \text{and} \quad \mathcal{R}(\mu_{\mathbf{W}}) := \int \|\cdot\|^2 d\mu_{\mathbf{W}} \quad \text{with} \quad \mu_{\mathbf{W}} = \frac{1}{m} \sum_{j=1}^m \delta_{\mathbf{w}_j}. \quad (\text{A.2})$$

Indeed, $\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) = \hat{y}_m(\mathbf{x}; \mathbf{W})$ and $\mathcal{R}(\mu_{\mathbf{W}}) = R(\mathbf{W})$, and this formulation allows extension to infinite-width networks by letting $\mu \in \mathcal{P}_2(\mathcal{W})$. Thus, we rewrite the population and the empirical risks in the space of measures as

$$\mathcal{J}_\lambda(\mu_{\mathbf{W}}) := J_\lambda(\mathbf{W}) \quad \text{and} \quad \hat{\mathcal{J}}_\lambda(\mu_{\mathbf{W}}) := \hat{J}_\lambda(\mathbf{W}),$$

and allow their domain to be all $\mu \in \mathcal{P}_2(\mathcal{W})$. We can equivalently state the interacting SDE system (A.1) as (see e.g. [16, Proposition 2.4])

$$d\mathbf{w}_j^t = -\nabla_{\mathbf{w}} \hat{\mathcal{J}}'_\lambda[\mu_{\mathbf{W}^t}](\mathbf{w}_j^t) + \sqrt{\frac{2}{\beta}}d\mathbf{B}_t^j \quad \text{for } 1 \leq j \leq m, \quad (\text{A.3})$$

where $\hat{\mathcal{J}}'_\lambda[\mu] \in L^2(\mathcal{W})$ denotes the first variation [42, Definition 7.12] of $\hat{\mathcal{J}}_\lambda(\mu)$.

As $m \rightarrow \infty$, the stochastic empirical measure $\mu_{\mathbf{W}^t}$ weakly converges to a deterministic measure μ_t for all fixed t , a phenomenon known as the *propagation of chaos* [46]. Furthermore, μ_t can be characterized as the law of the solution of the following SDE and non-linear Fokker-Planck equation

$$d\mathbf{w}^t = -\nabla_{\mathbf{w}} \hat{\mathcal{J}}'_\lambda[\mu_t](\mathbf{w}^t)dt + \sqrt{\frac{2}{\beta}}d\mathbf{B}_t \quad \text{and} \quad \partial_t \mu_t = \nabla \cdot (\mu_t \nabla \hat{\mathcal{J}}'_\lambda[\mu_t]) + \beta^{-1} \Delta \mu_t, \quad (\text{A.4})$$

where $\nabla \cdot$ and Δ are the divergence and Laplacian operators, respectively. Due to the existence of mean-field interactions, (A.4) is known as the *mean-field Langevin dynamics* (MFLD).

For a pair of probability measures $\mu \ll \nu$ both in $\mathcal{P}(\mathcal{W})$, we define the relative entropy $\mathcal{H}(\mu | \nu)$ and the relative Fisher information $\mathcal{I}(\mu | \nu)$ respectively as

$$\mathcal{H}(\mu | \nu) := \int_{\mathcal{W}} \ln \frac{d\mu}{d\nu} d\mu \quad \text{and} \quad \mathcal{I}(\mu | \nu) := \int_{\mathcal{W}} \left\| \nabla \ln \frac{d\mu}{d\nu} \right\|^2 d\mu. \quad (\text{A.5})$$

It is well-known at this point that μ_t in (A.4) can be interpreted as the Wasserstein gradient flow of the entropic regularized functional $\mathcal{F}_\beta(\mu) := \hat{\mathcal{J}}_\lambda(\mu) + \frac{1}{\beta}\mathcal{H}(\mu|\tau)$, where τ is the uniform measure on compact \mathcal{W} or the Lebesgue measure on a Euclidean space [3, 31, 49]. For this gradient flow to converge exponentially fast towards the minimizer $\mu_\beta^* := \arg \min_\mu \mathcal{F}_\beta(\mu)$, we require a *gradient domination* condition on μ_β^* in the space of probability measures, given as

$$\mathcal{H}(\mu|\mu_\beta^*) \leq \frac{C_{\text{LSI}}}{2}\mathcal{I}(\mu|\mu_\beta^*), \quad \forall \mu \in \mathcal{P}(\mathcal{W}), \quad (\text{A.6})$$

which is referred to as the log-Sobolev inequality (LSI). If the measure $d\nu_{\mu_t} \propto \exp(-\beta\mathcal{F}'[\mu_t])d\tau$ satisfies LSI with constant C_{LSI} for all $t \geq 0$, μ_t enjoys the following exponential convergence

$$\mathcal{F}_\beta(\mu_t) - \mathcal{F}_\beta(\mu_\beta^*) \leq e^{\frac{-2t}{\beta C_{\text{LSI}}}}(\mathcal{F}_\beta(\mu_0) - \mathcal{F}_\beta(\mu_\beta^*)); \quad (\text{A.7})$$

see e.g. [16, Theorem 3.2] and [38, Theorem 1].

Appendix B. Proofs of Section 3

Before presenting the layout of proofs, we introduce a useful reformulation of the objective $\mathcal{F}_{\beta,\lambda}(\mu)$. Recall that

$$\mathcal{F}_{\beta,\lambda}(\mu) = \hat{\mathcal{J}}_0(\mu) + \frac{\lambda}{2}\mathcal{R}(\mu) + \frac{1}{\beta}\mathcal{H}(\mu).$$

Let $\gamma \propto \exp\left(\frac{-\lambda\beta}{2}\|\mathbf{w}\|^2\right)$ be the centered Gaussian measure on \mathbb{R}^{2d+2} with variance $1/(\lambda\beta)$. Then, we can rewrite the above as

$$\mathcal{F}_{\beta,\lambda}(\mu) = \hat{\mathcal{J}}_0(\mu) + \frac{1}{\beta}\mathcal{H}(\mu|\gamma) - \frac{d}{2}\ln\left(\frac{\lambda\beta}{2}\right).$$

As a result, we can define

$$\tilde{\mathcal{F}}_{\beta,\lambda}(\mu) := \hat{\mathcal{J}}_0(\mu) + \frac{1}{\beta}\mathcal{H}(\mu|\gamma), \quad (\text{B.1})$$

which is non-negative and equivalent to \mathcal{F}_β up to an additive constant. Notice that

$$\mu_\beta^* := \arg \min_\mu \mathcal{F}_{\beta,\lambda}(\mu) = \arg \min_\mu \tilde{\mathcal{F}}_{\beta,\lambda}(\mu).$$

This reformulation, which was also used in [45], allows us to combine the effect of weight decay and entropic regularization into a single non-negative term $\mathcal{H}(\mu|\gamma)$. Furthermore, the simple density expression for the Gaussian measure γ allows us to achieve useful estimates for $\mathcal{H}(\mu|\gamma)$. In particular, as we will show below, it is possible to control $\mathcal{H}(\mu_\beta^*|\gamma)$ with effective dimension rather than ambient dimension, which leads to dependence on d_{eff} rather than d in our bounds.

We break down the proof of Theorem 2 into three steps:

1. In Section B.2 we show that there exists a measure $\mu^* \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ where $\hat{y}(\cdot; \mu^*)$ can approximate g on the training set with bounds on $\mathcal{R}(\mu^*)$. This construction provides upper bounds on $\hat{\mathcal{J}}_0(\mu_\beta^*)$ and $\mathcal{H}(\mu_\beta^*|\gamma)$.
2. In Section B.3, we perform a generalization analysis via Rademacher complexity tools given the bound on $\mathcal{H}(\mu_\beta^*|\gamma)$, leading to a bound on $\mathcal{J}_0(\mu_\beta^*)$.

3. Finally, in Section B.4, we estimate the LSI constant and constants related to smoothness/discretization along the trajectory, which imply that $\mathcal{F}_{\beta,\lambda}^m(\mu_l^m)$ converges to $\mathcal{F}_\beta(\mu_\beta^*)$, where $\mathcal{F}_{\beta,\lambda}^m$ is an adjusted objective over $\mathcal{P}(\mathbb{R}^{(2d+2)m})$ defined in (B.6). This bound implies the convergence of $\mathbb{E}_{\mathbf{W} \sim \mu_l^m}[J_0(\mathbf{W})]$ to $J_0(\mu_\beta^*)$, which was bounded in the previous step.

Before laying out these steps, in Section B.1, we will introduce the required concentration results. In the following, we will use the unregularized population $\mathcal{J}_0(\mu) := \mathbb{E}[\ell(\hat{y}(\mathbf{x}; \mu), y)]$ and empirical $\hat{\mathcal{J}}_0(\mu) = \mathbb{E}_{S_n}[\ell(\hat{y}(\mathbf{x}; \mu), y)]$ risks, and also consider the finite-width versions $J_0(\mathbf{W}) := \mathcal{J}_0(\mu_{\mathbf{W}})$ and $\hat{J}_0(\mathbf{W}) := \hat{\mathcal{J}}_0(\mu_{\mathbf{W}})$.

B.1. Concentration Bounds

We begin by specifying the definition of subGaussian and subexponential random variables in our setting.

Definition 5 [50] *A random variable x is σ -subGaussian if $\mathbb{E}[e^{\lambda(x - \mathbb{E}[x])}] \leq e^{\lambda^2 \sigma^2 / 2}$ for all $\lambda \in \mathbb{R}$, and is (ν, α) -subexponential if $\mathbb{E}[e^{\lambda(x - \mathbb{E}[x])}] \leq e^{\lambda^2 \nu^2 / 2}$ for all $|\lambda| \leq 1/\alpha$. If x is σ -subGaussian, then*

$$\mathbb{P}(x - \mathbb{E}[x] \geq t) \leq \exp\left(\frac{-t^2}{2\sigma^2}\right). \quad (\text{B.2})$$

If x is (ν, α) -subexponential, then

$$\mathbb{P}(x - \mathbb{E}[x] \geq t) \leq \exp\left(-\frac{1}{2} \min\left(\frac{t^2}{\nu^2}, \frac{t}{\alpha}\right)\right) \quad (\text{B.3})$$

Moreover, for centered random variables, let $|\cdot|_{\psi_2}$ and $|\cdot|_{\psi_1}$ denote the subGaussian and subexponential norm respectively [48, Definitions 2.5.6 and 2.7.5]. Then x is σ -subGaussian if and only if $\sigma \asymp |x - \mathbb{E}[x]|_{\psi_2}$, and is (ν, ν) -subexponential if and only if $\nu \asymp |x - \mathbb{E}[x]|_{\psi_1}$.

Next, we bound several quantities that appear in various parts of our proofs.

Lemma 6 *Under Assumption 1, for any $q > 0$ and all $1 \leq i \leq n$, with probability at least $1 - n^{-q}$,*

$$\left\| \mathbf{U}\mathbf{x}^{(i)} \right\| \leq r_x (1 + \sigma_u \sqrt{2(q+1) \ln n}) = \tilde{r}_x. \quad (\text{B.4})$$

Proof By subGaussianity of $\|\mathbf{U}\mathbf{x}\|$ from Assumption 1 and the subGaussian tail bound, with probability at least $1 - n^{-q-1}$

$$\begin{aligned} \left\| \mathbf{U}\mathbf{x}^{(i)} \right\| &\leq \mathbb{E}[\|\mathbf{U}\mathbf{x}\|] + \sigma_u r_x \sqrt{2(q+1) \ln n} \\ &= r_x + \sigma_u r_x \sqrt{2(q+1) \ln n}. \end{aligned}$$

The statement of lemma follows from a union bound over $1 \leq i \leq n$. ■

Lemma 7 *Under Assumption 1, we have $\mathbb{E}_{S_n} \left[\|\mathbf{x}\|^2 \right] \lesssim c_x^2$ with probability at least $1 - \exp(-\Omega(n))$.*

Proof By the triangle inequality,

$$\|\mathbf{x}\|_{\psi_2} \leq \|\mathbf{x}\| - \mathbb{E}[\|\mathbf{x}\|]_{\psi_2} + |\mathbb{E}[\|\mathbf{x}\|]|_{\psi_2} \lesssim \sigma_n \left\| \boldsymbol{\Sigma}^{1/2} \right\|_{\text{F}} + \text{tr}(\boldsymbol{\Sigma})^{1/2} \lesssim \text{tr}(\boldsymbol{\Sigma})^{1/2}.$$

Recall $c_x^2 := \text{tr}(\boldsymbol{\Sigma})$. Furthermore, by [48, Lemma 2.7.6] we have

$$\left| \|\mathbf{x}\|^2 \right|_{\psi_1} = \|\mathbf{x}\|_{\psi_2}^2 \lesssim c_x^2.$$

We arrive at a similar result for the centered random variable $\|\mathbf{x}\|^2 - \mathbb{E}[\|\mathbf{x}\|^2] = \|\mathbf{x}\|^2 - c_x^2$. We conclude the proof by the subexponential tail inequality,

$$\mathbb{P}\left(\mathbb{E}_{S_n} \left[\|\mathbf{x}\|^2 \right] - c_x^2 \geq tc_x^2\right) \leq \exp(-\min(t, t^2)\Omega(n)).$$

■

Lemma 8 *Under Assumption 1, we have $\mathbb{E}_{S_n} [y^2] \lesssim 1$ with probability at least $1 - 2\exp(-\Omega(n))$.*

Proof We have

$$|y|^2 \leq 3g(0)^2 + 3\mathcal{O}(1/r_x^2)\|\mathbf{U}\mathbf{x}\|^2 + 3\xi^2.$$

By a similar argument to Lemma 7 we have

$$\left| \|\mathbf{U}\mathbf{x}\|^2 \right|_{\psi_1} = \|\mathbf{U}\mathbf{x}\|_{\psi_2}^2 \leq 2\|\mathbf{U}\mathbf{x}\| - \mathbb{E}[\|\mathbf{U}\mathbf{x}\|]_{\psi_2}^2 + 2\mathbb{E}[\|\mathbf{U}\mathbf{x}\|^2] \lesssim (1 + \sigma_u^2)r_x^2,$$

since $\mathbb{E}[\|\mathbf{U}\mathbf{x}\|^2] = r_x^2$. As a result, by the subexponential tail bound,

$$\mathbb{E}_{S_n} \left[\|\mathbf{U}\mathbf{x}\|^2 \right] - \mathbb{E} \left[\|\mathbf{U}\mathbf{x}\|^2 \right] \lesssim (1 + \sigma_u^2)r_x^2 \lesssim r_x^2,$$

with probability at least $1 - \exp(-\Omega(n))$. Similarly, $|\xi^2|_{\psi_1} \leq |\xi|_{\psi_2}^2 \lesssim \varsigma^2$, therefore,

$$\mathbb{E}_{S_n} [\xi^2] - \mathbb{E}[\xi^2] \lesssim \varsigma^2 \lesssim 1,$$

with probability at least $1 - \exp(-\Omega(n))$. The statement of the lemma follows by a union bound. ■

Lemma 9 *Under Assumption 1, for any $q > 0$ and $n \gtrsim \frac{c_x^2}{\|\boldsymbol{\Sigma}\|} (1 + \sigma_n^2(q+1)\ln(n)) \ln(dn^q)$, with probability at least $1 - \mathcal{O}(n^{-q})$ we have $\|\mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top]\| \lesssim \|\boldsymbol{\Sigma}\|$. Further, if $q \geq 1$, then $\mathbb{E} \left[\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\|^{1/2} \right] \lesssim \|\boldsymbol{\Sigma}\|^{1/2}$.*

Proof First, note that by subGaussianity of $\|\mathbf{x}\|$, for every fixed i , we have with probability at least $1 - n^{-q-1}$,

$$\left| \mathbf{x}^{(i)} \right| - \mathbb{E}[\|\mathbf{x}\|] \leq \sigma_n \left\| \boldsymbol{\Sigma}^{1/2} \right\|_{\text{F}} \sqrt{2(q+1)\ln n}.$$

Since $\mathbb{E}[\|\mathbf{x}\|] \leq c_x$, via a union bound, with probability at least $1 - n^{-q}$,

$$\left\| \mathbf{x}^{(i)} \right\| \leq c_x + \sigma_n c_x \sqrt{2(q+1) \ln n} =: \tilde{c}_x.$$

Define the clipped version of \mathbf{x} via $\mathbf{x}_c = \mathbf{x}(1 \wedge \frac{\tilde{c}_x}{\|\mathbf{x}\|})$. Then, on the above event,

$$\mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] = \mathbb{E}_{S_n} [\mathbf{x}_c\mathbf{x}_c^\top].$$

Moreover,

$$\left\| \mathbb{E} [\mathbf{x}_c\mathbf{x}_c^\top] \right\| = \sup_{\|\mathbf{v}\| \leq 1} \mathbb{E} [\langle \mathbf{x}_c, \mathbf{v} \rangle^2] \leq \sup_{\|\mathbf{v}\| \leq 1} \mathbb{E} [\langle \mathbf{x}, \mathbf{v} \rangle^2] = \left\| \mathbb{E} [\mathbf{x}\mathbf{x}^\top] \right\|.$$

Finally, by the covariance estimation bound of [50, Corollary 6.20] for centered subGaussian random vectors and the condition on n given in the statement of the lemma,

$$\left\| \mathbb{E}_{S_n} [\mathbf{x}_c\mathbf{x}_c^\top] \right\| - \left\| \mathbb{E} [\mathbf{x}_c\mathbf{x}_c^\top] \right\| \lesssim \left\| \mathbb{E} [\mathbf{x}\mathbf{x}^\top] \right\|$$

with probability at least $1 - \mathcal{O}(n^{-q})$. Consequently, we have $\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\| \lesssim \|\boldsymbol{\Sigma}\|$ with probability at least $1 - \mathcal{O}(n^{-q})$.

For the second part of the lemma, let E denote the event on which the above $\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\| \lesssim \|\boldsymbol{\Sigma}\|$ holds. Then,

$$\begin{aligned} \mathbb{E} \left[\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\|^{1/2} \right] &= \mathbb{E} \left[\mathbb{1}(E) \left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\|^{1/2} \right] + \mathbb{E} \left[\mathbb{1}(E^C) \left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\|^{1/2} \right] \\ &\lesssim \|\boldsymbol{\Sigma}\|^{1/2} + \mathbb{P}(E^C)^{1/2} \mathbb{E} \left[\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\|^{1/2} \right] \\ &\lesssim \|\boldsymbol{\Sigma}\|^{1/2} + \mathcal{O}(n^{-q/2})c_x. \end{aligned}$$

Suppose $q \geq 1$. Then for $n \gtrsim c_x^2/\|\boldsymbol{\Sigma}\|$, we have $\mathbb{E} \left[\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\|^{1/2} \right] \lesssim \|\boldsymbol{\Sigma}\|^{1/2}$, which completes the proof. \blacksquare

We summarize the above results into a single event.

Lemma 10 *Suppose $n \gtrsim \frac{c_x^2}{\|\boldsymbol{\Sigma}\|} (1 + \sigma_n^2 (q+1) \ln(n)) \ln(dn^q)$. There exists an event \mathcal{E} such that $\mathbb{P}(\mathcal{E}) \geq 1 - \mathcal{O}(n^{-q})$, and on \mathcal{E} :*

1. $\|\mathbf{U}\mathbf{x}^{(i)}\| \leq \tilde{r}_x$ for all $1 \leq i \leq n$.
2. $\mathbb{E}_{S_n} [\|\mathbf{x}\|^2] \lesssim c_x^2$.
3. $\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\| \lesssim \|\boldsymbol{\Sigma}\|$.
4. $\mathbb{E} \left[\left\| \mathbb{E}_{S_n} [\mathbf{x}\mathbf{x}^\top] \right\|^{1/2} \right] \lesssim \|\boldsymbol{\Sigma}\|^{1/2}$.
5. $\mathbb{E}_{S_n} [y^2] \lesssim 1$.

We recall the variational lower bound for the KL divergence, which will be used at various stages of different proofs to relate certain expectations to the KL divergence.

Lemma 11 (Donsker-Varadhan Variational Formula for KL Divergence [26]) *Let μ and ν be probability measures on \mathcal{W} . Then,*

$$\mathcal{H}(\mu | \nu) = \sup_{f: \mathcal{W} \rightarrow \mathbb{R}} \int f d\mu - \ln \left(\int e^f d\nu \right).$$

Finally, we state the following lemma which will be useful in estimating smoothness constants in the convergence analysis.

Lemma 12 *Suppose $(z, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^d$ are drawn from a probability distribution \mathcal{D} . Then,*

$$\|\mathbb{E}_{\mathcal{D}}[z\mathbf{x}]\| \leq \sqrt{\mathbb{E}_{\mathcal{D}}[z^2] \|\mathbb{E}_{\mathcal{D}}[\mathbf{x}\mathbf{x}^{\top}]\|}.$$

Proof We have

$$\begin{aligned} \|\mathbb{E}_{\mathcal{D}}[z\mathbf{x}]\| &= \sup_{\|\mathbf{v}\| \leq 1} \langle \mathbf{v}, \mathbb{E}_{\mathcal{D}}[z\mathbf{x}] \rangle = \sup_{\|\mathbf{v}\| \leq 1} \mathbb{E}_{\mathcal{D}}[z \langle \mathbf{v}, \mathbf{x} \rangle] \\ &\leq \sup_{\|\mathbf{v}\| \leq 1} \sqrt{\mathbb{E}_{\mathcal{D}}[z^2] \mathbb{E}_{\mathcal{D}}[\langle \mathbf{v}, \mathbf{x} \rangle^2]} \quad (\text{Cauchy-Schwartz}) \\ &\leq \sqrt{\mathbb{E}_{\mathcal{D}}[z^2] \sup_{\|\mathbf{v}\| \leq 1} \langle \mathbf{v}, \mathbb{E}_{\mathcal{D}}[\mathbf{x}\mathbf{x}^{\top}]\mathbf{v} \rangle} \\ &= \sqrt{\mathbb{E}_{\mathcal{D}}[z^2] \|\mathbb{E}_{\mathcal{D}}[\mathbf{x}\mathbf{x}^{\top}]\|}. \end{aligned}$$

■

Notice that the distribution \mathcal{D} can be both the empirical as well as the population distribution.

B.2. Approximating the Target Function

We begin by stating the following approximation lemma which is the result of [7, Proposition 6] adapted to our setting.

Proposition 13 *Suppose $g : \mathbb{R}^k \rightarrow \mathbb{R}$ is L -Lipschitz and $|g(0)| = \mathcal{O}(L\tilde{r}_x)$. On the event of Lemma 10, there exists a measure $\mu \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ with $\mathcal{R}(\mu) \leq \Delta^2/\tilde{r}_x^2$ such that*

$$\max_i \left| g(\mathbf{U}\mathbf{x}^{(i)}) - \hat{y}(\mathbf{x}^{(i)}; \mu) \right| \leq C_k L \tilde{r}_x \left(\frac{\Delta}{L\tilde{r}_x} \right)^{\frac{-2}{k+1}} \ln \left(\frac{\Delta}{L\tilde{r}_x} \right) + \frac{\ln 4}{\kappa},$$

for all $\Delta \geq C_k$, where C_k is a constant depending only on k , provided that the hyperparameter ι satisfies $\iota \geq C_k L \tilde{r}_x \left(\frac{\Delta}{L\tilde{r}_x} \right)^{2k/(k+1)}$.

Proof Throughout the proof, we will use C_k to denote a constant that only depends on k , whose value may change across instantiations. Let $\mathbf{z} := \mathbf{U}\mathbf{x} \in \mathbb{R}^k$ and $\tilde{\mathbf{z}} := (\mathbf{z}^{\top}, \tilde{r}_x)^{\top} \in \mathbb{R}^{k+1}$. Recall that on the event of Lemma 6 we have $\|\mathbf{z}^{(i)}\| \leq \tilde{r}_x$ and $|g(\mathbf{z}^{(i)})| \lesssim L\tilde{r}_x$ for all $1 \leq i \leq n$. Let τ

denote the uniform probability measure on \mathbb{S}^k . By [7, Proposition 6], for all $\Delta \geq C_k$, there exists $p \in L^2(\tau)$ with $\|p\|_{L^2(\tau)} \leq \Delta$ such that

$$\max_i \left| g(\mathbf{z}^{(i)}) - \int_{\mathbb{S}^k} p(\mathbf{v}) \phi_\infty \left(\frac{1}{\tilde{r}_x} \langle \mathbf{v}, \tilde{\mathbf{z}}^{(i)} \rangle \right) d\tau(\mathbf{v}) \right| \leq C_k L \tilde{r}_x \left(\frac{\Delta}{L \tilde{r}_x} \right)^{\frac{-2}{k+1}} \ln \left(\frac{\Delta}{L \tilde{r}_x} \right).$$

In fact, we have a stronger guarantee on p . Specifically, $p(\mathbf{v})$ is given by

$$p(\mathbf{v}) = \sum_{j \geq 1} \lambda_j^{-1} r^j h_j(\mathbf{v}),$$

where $r \in (0, 1)$, $\lambda_j, h_j : \mathbb{S}^k \rightarrow \mathbb{R}$ are introduced by [7, Appendix D]. In particular,

$$h(\mathbf{v}) = g \left(\frac{\tilde{r}_x \mathbf{v}_{1:k}}{\mathbf{v}_{k+1}} \right) \mathbf{v}_{k+1},$$

with the spherical harmonics decomposition $h(\mathbf{v}) = \sum_{j \geq 0} h_j(\mathbf{v})$. It is shown in [7, Appendix D.2] that $\lambda_j \leq C_k j^{(k+1)/2}$, and one can prove through spherical harmonics calculations (omitted here for brevity) that $|h_j(\mathbf{v})| \leq C_k \sup_{\mathbf{v} \in \mathbb{S}^k} h(\mathbf{v}) j^{(k-1)/2} \leq C_k L \tilde{r}_x j^{(k-1)/2}$. As a result,

$$|p(\mathbf{v})| \leq \sum_{j \geq 0} \lambda_j^{-1} r^j |h_j(\mathbf{v})| \leq \sum_{j \geq 1} \lambda_j^{-1} r^j |h_j(\mathbf{v})| \leq C_k L \tilde{r}_x \sum_{j \geq 1} j^k r^j \leq \frac{C_k L \tilde{r}_x}{(1-r)^k}.$$

Using $1-r = \left(C_k L \tilde{r}_x / \Delta \right)^{2/(k+1)}$ as in [7, Appendix D.4] yields

$$|p(\mathbf{v})| \leq C_k L \tilde{r}_x \left(\frac{\Delta}{L \tilde{r}_x} \right)^{2k/(k+1)}.$$

Define $p_+(\mathbf{v}) := p(\mathbf{v}) \vee 0$ and $p_-(\mathbf{v}) := (-p(\mathbf{v})) \vee 0$. Then, by positive 1-homogeneity of ReLU,

$$\begin{aligned} \int_{\mathbb{S}^k} p(\mathbf{v}) \phi_\infty \left(\frac{1}{\tilde{r}_x} \langle \mathbf{v}, \tilde{\mathbf{z}} \rangle \right) d\tau(\mathbf{v}) &= \int_{\mathbb{S}^k} p_+(\mathbf{v}) \phi_\infty \left(\frac{1}{\tilde{r}_x} \langle \mathbf{v}, \tilde{\mathbf{z}} \rangle \right) d\tau(\mathbf{v}) - \int_{\mathbb{S}^k} p_-(\mathbf{v}) \phi_\infty \left(\frac{1}{\tilde{r}_x} \langle \mathbf{v}, \tilde{\mathbf{z}} \rangle \right) d\tau(\mathbf{v}) \\ &= \int_{\mathbb{S}^k} \phi_\infty \left(\frac{p_+(\mathbf{v})}{\tilde{r}_x} \langle \mathbf{v}, \tilde{\mathbf{z}} \rangle \right) d\tau(\mathbf{v}) - \int_{\mathbb{S}^k} \phi_\infty \left(\frac{p_-(\mathbf{v})}{\tilde{r}_x} \langle \mathbf{v}, \tilde{\mathbf{z}}^{(i)} \rangle \right) d\tau(\mathbf{v}) \\ &= \int_{\mathbb{R}^{k+1}} \phi_\infty(\langle \mathbf{v}, \tilde{\mathbf{z}} \rangle) d\tilde{\mu}_1(\mathbf{v}) - \int_{\mathbb{R}^{k+1}} \phi_\infty(\langle \mathbf{v}, \tilde{\mathbf{z}} \rangle) d\tilde{\mu}_2(\mathbf{v}) \\ &= \int_{\mathbb{R}^{d+1}} \phi_\infty(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) d\mu_1(\mathbf{w}) - \int_{\mathbb{R}^{d+1}} \phi_\infty(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) d\mu_2(\mathbf{w}), \end{aligned}$$

where $\tilde{\mu}_1 := \frac{(\cdot) p_+(\cdot)}{\tilde{r}_x} \# \tau$ and $\tilde{\mu}_2 := \frac{(\cdot) p_-(\cdot)}{\tilde{r}_x} \# \tau$ are the corresponding pushforward measures, $\mu_1 = T_U \# \tilde{\mu}_1$ and $\mu_2 = T_U \# \tilde{\mu}_2$, where $T_U(\mathbf{v}) = (\mathbf{U}^\top \mathbf{v}_k, v_{k+1})^\top \in \mathbb{R}^{d+1}$ for $\mathbf{v} = (\mathbf{v}_k^\top, v_{k+1})^\top \in \mathbb{R}^{k+1}$. In other words, $\mathbf{w} \sim \mu_1$ is generated by sampling $\mathbf{v} \sim \tilde{\mu}_1$ and letting $\mathbf{w} = (\mathbf{U}^\top \mathbf{v}_k, v_{k+1})^\top$, with a similar procedure for $\mathbf{w} \sim \mu_2$. Furthermore,

$$\begin{aligned} \mathcal{R}(\mu) &= \int_{\mathbb{R}^{d+1}} \|\mathbf{w}\|^2 d\mu_1(\mathbf{w}) + \int_{\mathbb{R}^{d+1}} \|\mathbf{w}\|^2 d\mu_2(\mathbf{w}) = \int_{\mathbb{R}^{k+1}} \|\mathbf{v}\|^2 d\tilde{\mu}_1(\mathbf{v}) + \int_{\mathbb{R}^{k+1}} \|\mathbf{v}\|^2 d\tilde{\mu}_2(\mathbf{v}) \\ &= \int_{\mathbb{S}^k} \frac{p(\mathbf{v})^2}{\tilde{r}_x^2} d\tau(\mathbf{v}) \leq \frac{\Delta^2}{\tilde{r}_x^2}. \end{aligned}$$

The last step is to replace ϕ_∞ with $\phi_{\kappa,\iota}$. Note that for all i , and almost surely over $\mathbf{w} \sim \mu_1$, we have $\left| \langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle \right| \leq p_+(\mathbf{v}) \leq C_k L \tilde{r}_x \left(\frac{\Delta}{L \tilde{r}_x} \right)^{2k/(k+1)}$, with a similar bound holding for $\mathbf{w} \sim \mu_2$. As a result, by choosing $\iota \geq C_k L \tilde{r}_x \left(\frac{\Delta}{L \tilde{r}_x} \right)^{2k/(k+1)}$, we have $\phi_{\kappa,\iota}(\langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle) = \phi_\infty(\langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle)$ for all i and almost surely over $\mathbf{w} \sim \mu_1$ and $\mathbf{w} \sim \mu_2$. By the triangle inequality, we have

$$\begin{aligned} \left| g(\mathbf{U}\mathbf{x}^{(i)}) - \hat{y}(\mathbf{x}^{(i)}; \mu) \right| &\leq \left| \int \phi_{\kappa,\iota}(\langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle) - \phi_\infty(\langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle) d\mu_1(\mathbf{w}) \right| \\ &\quad + \left| \int \phi_{\kappa,\iota}(\langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle) - \phi_\infty(\langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle) d\mu_2(\mathbf{w}) \right| \\ &\quad + \left| g(\mathbf{U}\mathbf{x}^{(i)}) - \int \phi_\infty(\langle \mathbf{w}, \tilde{\mathbf{x}}^{(i)} \rangle) (d\mu_1(\mathbf{w}) - d\mu_2(\mathbf{w})) \right| \\ &\leq \frac{2 \ln 2}{\kappa} + C_k L \tilde{r}_x \left(\frac{\Delta}{L \tilde{r}_x} \right)^{\frac{-2}{k+1}} \ln \left(\frac{\Delta}{L \tilde{r}_x} \right), \end{aligned}$$

which completes the proof. \blacksquare

Next, we control the effect of entropic regularization on the minimum of $\tilde{\mathcal{F}}_{\beta,\lambda}$ via the following lemma.

Lemma 14 *Suppose ρ is C_ρ Lipschitz. For every $\mu^* \in \mathcal{P}(\mathbb{R}^{2d+2})$, we have*

$$\min_{\mu \in \mathcal{P}^{\text{ac}}(\mathbb{R}^{2d+2})} \tilde{\mathcal{F}}_{\beta,\lambda}(\mu) \leq \hat{\mathcal{J}}_0(\mu^*) + \frac{\lambda}{2} \mathcal{R}(\mu^*) + \frac{2\sqrt{2}C_\rho}{\sqrt{\pi\lambda\beta}} \mathbb{E}_{S_n}[\|\tilde{\mathbf{x}}\|].$$

Proof We will smooth μ^* by convolving it with γ , i.e. we consider $\mu = \mu^* * \gamma$. Let $\mathbf{u} \sim \gamma$ independent of $\mathbf{w} \sim \mu^*$ and denote $\mathbf{u} = (\mathbf{u}_1^\top, \mathbf{u}_2^\top)^\top$ with $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^{d+1}$. We first bound $\hat{\mathcal{J}}_0(\mu^* * \gamma)$. Using the Lipschitzness of the loss and of $\phi_{\kappa,\iota}$, we have

$$\begin{aligned} \hat{\mathcal{J}}_0(\mu^* * \gamma) - \hat{\mathcal{J}}_0(\mu^*) &= \mathbb{E}_{S_n} \left[\ell \left(\int \Psi(\mathbf{x}; \mathbf{w}) d(\mu^* * \gamma)(\mathbf{w}) - y \right) - \ell \left(\int \Psi(\mathbf{x}; \mathbf{w}) d\mu^*(\mathbf{w}) - y \right) \right] \\ &\leq C_\rho \mathbb{E}_{S_n} \left[\left| \int \Psi(\mathbf{x}; \mathbf{w}) d(\mu^* * \gamma)(\mathbf{w}) - \int \Psi(\mathbf{x}; \mathbf{w}) d\mu^*(\mathbf{w}) \right| \right] \\ &= C_\rho \mathbb{E}_{S_n} \left[\left| \int (\mathbb{E}_{\mathbf{u}}[\Psi(\mathbf{x}; \mathbf{w} + \mathbf{u})] - \Psi(\mathbf{x}; \mathbf{w})) d\mu^*(\mathbf{w}) \right| \right] \\ &\leq C_\rho \mathbb{E}_{S_n} \left[\int \mathbb{E}_{\mathbf{u}} [|\phi_{\kappa,\iota}(\langle \boldsymbol{\omega}_1 + \mathbf{u}_1, \tilde{\mathbf{x}} \rangle) - \phi_{\kappa,\iota}(\langle \boldsymbol{\omega}_1, \tilde{\mathbf{x}} \rangle)|] d\mu^*(\mathbf{w}) \right] \\ &\quad + C_\rho \mathbb{E}_{S_n} \left[\int \mathbb{E}_{\mathbf{u}} [|\phi_{\kappa,\iota}(\langle \boldsymbol{\omega}_2 + \mathbf{u}_2, \tilde{\mathbf{x}} \rangle) - \phi_{\kappa,\iota}(\langle \boldsymbol{\omega}_2, \tilde{\mathbf{x}} \rangle)|] d\mu^*(\mathbf{w}) \right] \\ &\leq C_\rho \mathbb{E}_{S_n} \left[\int \{ \mathbb{E}_{\mathbf{u}_1} [|\langle \mathbf{u}_1, \tilde{\mathbf{x}} \rangle|] + \mathbb{E}_{\mathbf{u}_2} [|\langle \mathbf{u}_2, \tilde{\mathbf{x}} \rangle|] \} d\mu^*(\boldsymbol{\omega}) \right] \\ &= \frac{2\sqrt{2}C_\rho}{\sqrt{\pi\lambda\beta}} \mathbb{E}_{S_n}[\|\tilde{\mathbf{x}}\|]. \end{aligned}$$

Next, we bound the KL divergence via its convexity in the first argument,

$$\mathcal{H}(\mu^* * \gamma | \gamma) = \mathcal{H} \left(\int \gamma(\cdot - \mathbf{w}') d\mu^*(\mathbf{w}') | \gamma \right) \leq \int \mathcal{H}(\gamma(\cdot - \mathbf{w}') | \gamma(\cdot)) d\mu^*(\mathbf{w}').$$

Furthermore,

$$\mathcal{H}(\gamma(\cdot - \mathbf{w}') | \gamma(\cdot)) = \int \frac{\lambda\beta}{2} (-\|\mathbf{w} - \mathbf{w}'\|^2 + \|\mathbf{w}\|^2) \gamma(d\mathbf{w} - \mathbf{w}') = \frac{\lambda\beta\|\mathbf{w}'\|^2}{2}.$$

Consequently,

$$\mathcal{H}(\mu^* * \gamma | \gamma) \leq \frac{\lambda\beta}{2} \mathcal{R}(\mu^*),$$

which finishes the proof. \blacksquare

Combining above results, we have the following statement.

Corollary 15 *Suppose the event of Lemma 10 holds, ρ is C_ρ Lipschitz, and $\lambda \lesssim 1$. Then,*

$$\min_{\mu \in \mathcal{P}^{\text{ac}}(\mathbb{R}^{2d+2})} \tilde{\mathcal{F}}_{\beta, \lambda}(\mu) - \mathbb{E}_{S_n}[\rho(\xi)] \lesssim C_\rho \frac{\tilde{r}_x}{r_x} \left(\frac{r_x \Delta}{\tilde{r}_x} \right)^{\frac{-2}{k+1}} \ln \left(\frac{r_x \Delta}{\tilde{r}_x} \right) + \frac{C_\rho}{\kappa} + \frac{\lambda \Delta^2}{\tilde{r}_x^2} + \frac{C_\rho(c_x + \tilde{r}_x)}{\sqrt{\lambda\beta}},$$

for all $\Delta \geq C_k$, provided that $\iota \geq C_k \Delta^{2k/(k+1)} (r_x/\tilde{r}_x)^{(k-1)/(k+1)}$.

Proof We will use Lemma 14 with $\mu^* \in \mathcal{P}(\mathbb{R}^{2d+2})$ constructed in Proposition 13. Then, for all $\Delta \geq C_k$,

$$\begin{aligned} \hat{\mathcal{J}}_0(\mu^*) &= \mathbb{E}[\rho(\hat{y}(\mathbf{x}; \mu^*) - y)] \\ &= \mathbb{E}_{S_n}[\rho(\hat{y}(\mathbf{x}; \mu^*) - g(\mathbf{U}\mathbf{x}) - \xi)] \\ &\leq \mathbb{E}_{S_n}[\rho(\xi)] + C_\rho \mathbb{E}_{S_n}[|\hat{y}(\mathbf{x}; \mu^*) - g(\mathbf{U}\mathbf{x})|] \\ &\leq \mathbb{E}_{S_n}[\rho(\xi)] + C_k C_\rho \frac{\tilde{r}_x}{r_x} \left(\frac{r_x \Delta}{\tilde{r}_x} \right)^{\frac{-2}{k+1}} \ln \left(\frac{r_x \Delta}{\tilde{r}_x} \right) + \frac{C_\rho \ln 4}{\kappa}. \end{aligned}$$

Furthermore, Proposition 13 guarantees $\mathcal{R}(\mu^*) \leq \Delta^2/\tilde{r}_x^2$. Combining these bounds with Lemma 14 completes the proof. \blacksquare

B.3. Generalization Analysis

Let

$$\mu_\beta^* := \arg \min_{\mu \in \mathcal{P}_2^{\text{ac}}(\mathbb{R}^{2d+2})} \mathcal{F}_\beta(\mu) = \arg \min_{\mu \in \mathcal{P}_2^{\text{ac}}(\mathbb{R}^{2d+2})} \tilde{\mathcal{F}}_\beta(\mu).$$

Corollary 15 gives an upper bound on $\hat{\mathcal{J}}_0(\mu^*)$. In this section, we transfer the bound to $\mathcal{J}_0(\mu^*)$ via a Rademacher complexity analysis. Since Corollary 15 implies a bound on $\mathcal{H}(\mu | \gamma)$, we will control the following quantity,

$$\sup_{\mu: \mathcal{H}(\mu | \gamma) \leq \Delta^2} \mathcal{J}_0(\mu) - \hat{\mathcal{J}}_0(\mu).$$

To be able to provide guarantees with high probability, we will prove uniform convergence over a truncated version of the risk instead, given by

$$\sup_{\mu: \mathcal{H}(\mu | \gamma) \leq \Delta^2} \mathcal{J}_0^\times(\mu) - \hat{\mathcal{J}}_0^\times(\mu),$$

where

$$\mathcal{J}_0^\varkappa(\mu) := \mathbb{E}[\rho_\varkappa(\hat{y}(\mathbf{x}; \mu) - y)], \quad \hat{\mathcal{J}}_0^\varkappa(\mu) := \mathbb{E}_{S_n}[\rho_\varkappa(\hat{y}(\mathbf{x}; \mu) - y)],$$

and $\rho_\varkappa(\cdot) := \rho(\cdot) \wedge \varkappa$. We will later specify the choice of \varkappa .

We are now ready to present the Rademacher complexity bound.

Lemma 16 ([14, Lemma 5.5], [45, Lemma 1]) *Suppose ρ is either a C_ρ -Lipschitz loss or the squared error loss. Let $\vartheta := \sqrt{2\varkappa}$ for the squared error loss and C_ρ for the Lipschitz loss. Recall $\gamma = \mathcal{N}(0, \frac{\mathbf{I}_{d+1}}{\lambda\beta})$. Then,*

$$\mathbb{E} \left[\sup_{\{\mu \in \mathcal{P}^{\text{ac}}(\mathbb{R}^{2d+2}) : \mathcal{H}(\mu | \gamma) \leq M\}} \mathcal{J}_0^\varkappa(\mu) - \hat{\mathcal{J}}_0^\varkappa(\mu) \right] \leq 4\vartheta \iota \sqrt{\frac{2M}{n}}.$$

Proof We repeat the proof here for the reader's convenience. Let $(\xi_i)_{i=1}^n$ denote i.i.d. Rademacher random variables. Notice that for the squared error loss, ρ_\varkappa is $\sqrt{2\varkappa}$ Lipschitz. Then, by a standard symmetrization argument and Talagrand's contraction lemma, we have

$$\begin{aligned} \mathbb{E} \left[\sup_{\mu : \mathcal{H}(\mu | \gamma) \leq M} \mathcal{J}_0(\mu) - \hat{\mathcal{J}}_0(\mu) \right] &\leq 2 \mathbb{E} \left[\sup_{\mu : \mathcal{H}(\mu | \gamma) \leq M} \frac{1}{n} \sum_{i=1}^n \xi_i \rho(\hat{y}(\mathbf{x}^{(i)}; \mu) - y) \right] \\ &\leq 2\vartheta \mathbb{E} \left[\sup_{\mu : \mathcal{H}(\mu | \gamma) \leq M} \frac{1}{n} \sum_{i=1}^n \xi_i \hat{y}(\mathbf{x}^{(i)}; \mu) \right] \end{aligned}$$

Next, we proceed to bound the Rademacher complexity. Specifically,

$$\begin{aligned} \mathbb{E}_\xi \left[\sup_{\mu : \mathcal{H}(\mu | \gamma) \leq M} \frac{1}{n} \sum_{i=1}^n \xi_i \int \Psi(\mathbf{x}^{(i)}; \mathbf{w}) d\mu(\mathbf{w}) \right] &= \mathbb{E}_\xi \left[\frac{1}{\alpha} \sup_{\mu : \mathcal{H}(\mu | \gamma) \leq M} \int \frac{\alpha}{n} \sum_{i=1}^n \xi_i \Psi(\mathbf{x}^{(i)}; \mathbf{w}) d\mu(\mathbf{w}) \right] \\ &\leq \frac{M}{\alpha} + \frac{1}{\alpha} \mathbb{E}_\xi \left[\ln \int \exp \left(\frac{\alpha}{n} \sum_{i=1}^n \xi_i \Psi(\mathbf{x}^{(i)}; \mathbf{w}) \right) d\gamma(\mathbf{w}) \right] \\ &\leq \frac{M}{\alpha} + \frac{1}{\alpha} \ln \int \mathbb{E}_\xi \left[\exp \left(\frac{\alpha}{n} \sum_{i=1}^n \xi_i \Psi(\mathbf{x}^{(i)}; \mathbf{w}) \right) \right] d\gamma(\mathbf{w}), \end{aligned}$$

where the first inequality follows from the KL divergence lower bound of Lemma 11. Additionally, by sub-Gaussianity and independence of (ξ_i) and Lipschitzness of $\phi_{\kappa, \iota}$, we have

$$\begin{aligned} \mathbb{E}_\xi \left[\exp \left(\frac{\alpha}{n} \sum_{i=1}^n \xi_i \Psi(\mathbf{x}^{(i)}; \mathbf{w}) \right) \right] &\leq \exp \left(\frac{\alpha^2}{2n^2} \sum_{i=1}^n \Psi(\mathbf{x}^{(i)}; \mathbf{w})^2 \right) \\ &\leq \exp \left(\frac{2\alpha^2 \iota^2}{n} \right) \end{aligned}$$

Plugging this back into our original bound, we obtain

$$\mathbb{E}_\xi \left[\sup_{\mu : \mathcal{H}(\mu | \gamma) \leq M} \frac{1}{n} \sum_{i=1}^n \xi_i \hat{y}(\mathbf{x}; \mu) \right] \leq \frac{M}{\alpha} + \frac{2\alpha \iota^2}{n}.$$

Choosing $\alpha = \sqrt{\frac{Mn}{2\iota^2}}$, we obtain

$$\mathbb{E}_{\xi} \left[\sup_{\mu: \mathcal{H}(\mu | \gamma) \leq M} \frac{1}{n} \sum_{i=1}^n \xi_i \hat{y}(\mathbf{x}; \mu) \right] \leq 2\iota \sqrt{\frac{2M}{n}},$$

which completes the proof. \blacksquare

We can convert the above bound in expectation to a high-probability bound as follows.

Lemma 17 *In the setting of Lemma 16, for any $\delta > 0$, we have*

$$\sup_{\mu \in \mathcal{P}^{\text{ac}}(\mathbb{R}^{2d+2}); \mathcal{H}(\mu | \gamma) \leq M} \mathcal{J}_0^{\varkappa}(\mu) - \hat{\mathcal{J}}_0^{\varkappa}(\mu) \lesssim \vartheta \iota \sqrt{\frac{M}{n}} + \varkappa \sqrt{\frac{\ln(1/\delta)}{n}},$$

with probability at least $1 - \delta$.

Proof As the truncated loss is bounded by \varkappa , the result is an immediate consequence of McDiarmid's inequality. \blacksquare

Next, we control the effect of truncation by bounding $\mathcal{J}_0(\mu)$ via $\mathcal{J}_0^{\varkappa}(\mu)$, which is achieved via the following lemma.

Lemma 18 *Suppose $\mathcal{H}(\mu | \gamma) \leq M$. Then,*

$$\mathcal{J}_0(\mu) - \mathcal{J}_0^{\varkappa}(\mu) \lesssim \left(\iota + \mathbb{E}[y^2]^{1/2} \right) e^{-\Omega(\varkappa^2)}.$$

Proof Notice that since the loss is C_ρ -Lipschitz and $\rho(0) = 0$, we have $|\rho(\hat{y} - y)| \leq C_\rho |\hat{y} - y|$. Recall that we use L for the Lipschitz constant of g , and $|\hat{y}(\mathbf{x}; \mu)| \leq 2\iota$. Then,

$$\begin{aligned} \mathcal{J}_0(\mu) - \mathcal{J}_0^{\varkappa}(\mu) &\leq \mathbb{E}[\mathbb{1}(\rho(\hat{y}(\mathbf{x}; \mu) - y) \geq \varkappa) \rho(\hat{y}(\mathbf{x}; \mu) - y)] \\ &\leq C_\rho \mathbb{P}(\rho(\hat{y}(\mathbf{x}; \mu) - y) \geq \varkappa)^{1/2} \mathbb{E}[(\hat{y}(\mathbf{x}; \mu) - y)^2]^{1/2} \\ &\leq C_\rho \mathbb{P}(2\iota + |y| \geq \varkappa/C_\rho)^{1/2} \left(\mathbb{E}[\hat{y}(\mathbf{x}; \mu)^2]^{1/2} + \mathbb{E}[y^2]^{1/2} \right) \\ &\leq C_\rho \mathbb{P}(2\iota + |g(0)| + L\|\mathbf{U}\mathbf{x}\| + \xi \geq \varkappa/C_\rho)^{1/2} \left(\mathbb{E}[\hat{y}(\mathbf{x}; \mu)^2]^{1/2} + \mathbb{E}[y^2]^{1/2} \right). \end{aligned}$$

Let $\varkappa/C_\rho \geq 4\iota + 2|g(0)| + 2Lr_x$, and recall that $L = \mathcal{O}(1/r_x)$. Then, by a subGaussian concentration bound, we have

$$\mathbb{P}(2\iota + |g(0)| + L\|\mathbf{U}\mathbf{x}\| + \xi \geq \varkappa/C_\rho)^{1/2} \leq e^{-\Omega\left(\frac{\varkappa^2}{\sigma_u^2 C_\rho^2}\right)}.$$

We conclude the proof by remarking that by our assumptions, σ_u and C_ρ are absolute constants. \blacksquare

Finally, we combine the steps above to give an upper bound on $\mathcal{J}_0(\mu_\beta^*)$, stated in the following lemma.

Lemma 19 Suppose $\lambda = \tilde{\lambda} r_x^2$ and $\beta = \frac{d_{\text{eff}} + \tilde{r}_x^2 / r_x^2}{\varepsilon^2 \tilde{\lambda}}$ for $\varepsilon, \tilde{\lambda} \lesssim 1$. Let $\tilde{\varepsilon} := \tilde{\mathcal{O}}(\tilde{\lambda}^{\frac{1}{k+2}}) + \varepsilon + \kappa^{-1}$. Suppose $n \gtrsim \frac{(d_{\text{eff}} + \tilde{r}_x^2 / r_x^2) \iota^2}{\lambda \varepsilon^4}$ and $\iota \gtrsim \tilde{\lambda}^{-\frac{k}{k+2}} (\tilde{r}_x / r_x)^{\frac{2k^2 + 4k + 2}{k+2}}$. Then,

$$\mathcal{J}_0(\mu_\beta^*) - \mathbb{E}[\rho(\xi)] \lesssim \tilde{\varepsilon}, \quad \text{and} \quad \beta^{-1} \mathcal{H}(\mu_\beta^* | \gamma) \lesssim \mathbb{E}[\rho(\xi)] + \tilde{\varepsilon} \lesssim 1.$$

Proof By Corollary 15 and a standard concentration bound on $\mathbb{E}_{S_n}[\rho(\xi)]$ with sufficiently large n to induce negligible error in comparison with the rest of the terms in the corollary, we have

$$\hat{\mathcal{J}}_0(\mu_\beta^*) + \beta^{-1} \mathcal{H}(\mu_\beta^* | \gamma) - \mathbb{E}[\rho(\xi)] \lesssim \frac{\tilde{r}_x}{r_x} \left(\frac{r_x \Delta}{\tilde{r}_x} \right)^{\frac{-2}{k+1}} \ln \left(\frac{r_x \Delta}{\tilde{r}_x} \right) + \frac{\lambda \Delta^2}{\tilde{r}_x^2} + \frac{(c_x + \tilde{r}_x)}{\sqrt{\lambda \beta}} + \frac{1}{\kappa}.$$

By choosing

$$\Delta = \left(\frac{r_x^2}{\lambda} \right)^{\frac{1}{2} \cdot \frac{k+1}{k+2}} \left(\frac{\tilde{r}_x}{r_x} \right)^{\frac{1}{2} \cdot \frac{3k+5}{k+2}},$$

and assuming $c_x \gtrsim \tilde{r}_x$,

$$\beta^{-1} \mathcal{H}(\mu_\beta^* | \gamma) \lesssim \mathbb{E}[\rho(\xi)] + \left(\frac{\lambda}{r_x^2} \right)^{\frac{1}{k+2}} \left(\frac{\tilde{r}_x}{r_x} \right)^{\frac{k+1}{k+2}} \ln \left(\frac{\tilde{r}_x r_x}{\lambda} \right) + \frac{c_x}{\sqrt{\lambda \beta}} + \frac{1}{\kappa}.$$

Note that the above choice on Δ translates to a lower bound on ι in Corollary 15, given by

$$\iota \gtrsim$$

By choosing $\lambda = \tilde{\lambda} r_x^2$ and using the fact that $\tilde{r}_x \leq \tilde{\mathcal{O}}(r_x)$ and $\beta = \frac{c_x^2}{r_x^2 \tilde{\lambda} \varepsilon^2}$, we have the simplification,

$$\beta^{-1} \mathcal{H}(\mu_\beta^* | \gamma) \lesssim \mathbb{E}[\rho(\xi)] + \tilde{\mathcal{O}}(\tilde{\lambda}^{\frac{1}{k+2}}) + \varepsilon + \frac{1}{\kappa} \lesssim 1,$$

and,

$$\hat{\mathcal{J}}_0(\mu_\beta^*) - \mathbb{E}[\rho(\xi)] \lesssim \tilde{\mathcal{O}}(\tilde{\lambda}^{\frac{1}{k+2}}) + \varepsilon + \frac{1}{\kappa} =: \tilde{\varepsilon}.$$

Note that $\hat{\mathcal{J}}_0^\varkappa(\mu_\beta^*) \leq \hat{\mathcal{J}}_0(\mu_\beta^*)$. Using the generalization bound of Lemma 17 with the choice of $\delta = n^{-q}$ for some constant $q > 0$, we have with probability $1 - \mathcal{O}(n^{-q})$,

$$\begin{aligned} \mathcal{J}_0^\varkappa(\mu_\beta^*) - \hat{\mathcal{J}}_0^\varkappa(\mu_\beta^*) &\lesssim \iota \sqrt{\frac{\beta}{n}} + \varkappa \sqrt{\frac{\ln n}{n}} \\ &\lesssim \iota \sqrt{\frac{d_{\text{eff}}}{n \tilde{\lambda} \varepsilon^2}} + \varkappa \sqrt{\frac{\ln n}{n}}. \end{aligned} \tag{B.5}$$

Furthermore, by Lemma 18 we have

$$\mathcal{J}_0(\mu_\beta^*) - \mathcal{J}_0^\varkappa(\mu_\beta^*) \lesssim \iota e^{-\Omega(\varkappa^2)}.$$

Combining the above with (B.5) and choosing on $\varkappa \asymp \sqrt{\ln n}$, we have

$$\mathcal{J}_0(\mu_\beta^*) - \mathbb{E}[\rho(\xi)] \lesssim \tilde{\varepsilon} + \iota \sqrt{\frac{d_{\text{eff}}}{n \tilde{\lambda} \varepsilon^2}} + \sqrt{\frac{\ln^2 n}{n}},$$

which holds with probability at least $1 - \mathcal{O}(n^{-q})$ over the randomness of S_n . \blacksquare

B.4. Convergence Analysis

So far, our analysis has only proved properties of μ_β^* . In this section, we relate these properties to μ_l^m via propagation of chaos. In particular, [44] showed that for $\mathbf{W} \sim \mu_l^m$, $\hat{y}(\mathbf{x}; \mu_l^m)$ converges to $\hat{y}(\mathbf{x}; \mu_\beta^*)$ in a suitable sense characterized shortly, as long as the objective over μ_l^m converges to $\mathcal{F}_{\beta,\lambda}(\mu_\beta^*)$. Notice that μ_l^m is a measure on $\mathcal{P}(\mathbb{R}^{(2d+2)m})$ instead of $\mathcal{P}(\mathbb{R}^{2d+2})$. Thus, we need to adjust the definition of objective by defining the following

$$\mathcal{F}_{\beta,\lambda}^m(\mu^m) := \mathbb{E}_{\mathbf{W} \sim \mu^m} \left[\hat{J}_0(\mathbf{W}) + \frac{\lambda}{2} R(\mathbf{W}) \right] + \frac{1}{m\beta} \mathcal{H}(\mu^m). \quad (\text{B.6})$$

We can use the same reformulation introduced earlier in (B.1) to define

$$\tilde{\mathcal{F}}_{\beta,\lambda}^m(\mu^m) := \mathbb{E}_{\mathbf{W} \sim \mu^m} \left[\hat{J}_0(\mathbf{W}) \right] + \frac{1}{m\beta} \mathcal{H}(\mu^m | \gamma^{\otimes m}), \quad (\text{B.7})$$

which is equivalent to $\mathcal{F}_{\beta,\lambda}^m$ up to an additive constant. With these definitions, we can now control $\mathbb{E}_{\mathbf{W} \sim \mu_l^m} [J_0(\mu_l^m)]$ via $\mathcal{J}_0(\mu_\beta^*)$. The following lemma is based on [44, Lemma 4], with a more careful analysis to obtain sharper constants.

Lemma 20 *Let $\bar{r}_x := \|\Sigma\|^{1/2} \vee \tilde{r}_x$, and suppose ρ is $C_\rho \lesssim 1$ -Lipschitz. Then,*

$$\mathbb{E}_{\mathbf{W} \sim \mu_l^m} [J_0(\mathbf{W})] - \mathcal{J}_0(\mu_\beta^*) \lesssim \sqrt{\frac{\bar{r}_x^2 W_2^2(\mu_l^m, \mu_\beta^{*\otimes m}) + \iota^2}{m}}. \quad (\text{B.8})$$

In particular, combined with [44, Lemma 3], the above implies

$$\mathbb{E}_{\mathbf{W} \sim \mu_l^m} [J_0(\mathbf{W})] - \mathcal{J}_0(\mu_\beta^*) \lesssim \sqrt{\frac{\bar{r}_x^2 \beta C_{\text{LSI}}}{m} (\tilde{\mathcal{F}}_{\beta,\lambda}^m(\mu_l^m) - \tilde{F}_{\beta,\lambda}(\mu_\beta^*))} + \frac{\iota^2}{m}. \quad (\text{B.9})$$

Proof Notice that

$$\begin{aligned} \mathbb{E}_{\mathbf{W} \sim \mu_l^m} [J_0(\mathbf{W})] &= \mathbb{E}_{\mathbf{W}} [\mathbb{E}_{\mathbf{x}} [\rho(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_\beta^*) + \hat{y}(\mathbf{x}; \mu_\beta^*) - y)]] \\ &\leq \mathbb{E}_{\mathbf{x}} [\rho(\hat{y}(\mathbf{x}; \mu_\beta^*) - y)] + C_\rho \mathbb{E}_{\mathbf{W}} [\mathbb{E}_{\mathbf{x}} [|\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_\beta^*)|]] \\ &\leq \mathcal{J}_0(\mu_\beta^*) + C_\rho \sqrt{\mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\mathbf{W}} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_\beta^*))^2]]} \end{aligned}$$

Suppose $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m) \sim \mu_l^m$ and $\mathbf{W}' = (\mathbf{w}'_1, \dots, \mathbf{w}'_m) \sim \mu_\beta^{*\otimes m}$. Let Γ denote the optimal W_2 coupling between \mathbf{W} and \mathbf{W}' , and assume $\mathbf{W}, \mathbf{W}' \sim \Gamma$. Then,

$$\begin{aligned} \mathbb{E}_{\mathbf{W}} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_\beta^*))^2] &= \mathbb{E}_{\mathbf{W}, \mathbf{W}'} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_{\mathbf{W}'}) + \hat{y}(\mathbf{x}; \mu_{\mathbf{W}'}) - \hat{y}(\mathbf{x}; \mu_\beta^*))^2] \\ &\leq 2 \mathbb{E}_{\mathbf{W}, \mathbf{W}'} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_{\mathbf{W}'})^2] + 2 \mathbb{E}_{\mathbf{W}'} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}'}) - \hat{y}(\mathbf{x}; \mu_\beta^*))^2] \end{aligned}$$

Moreover, by Jensen's inequality,

$$\begin{aligned} \mathbb{E}_{\mathbf{W}, \mathbf{W}'} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_{\mathbf{W}'})^2] &\leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathbf{W}, \mathbf{W}'} [(\Psi(\mathbf{x}; \mathbf{w}_i) - \Psi(\mathbf{x}; \mathbf{w}'_i))^2] \\ &\leq \frac{2}{m} \sum_{i=1}^m \mathbb{E}_{\mathbf{W}, \mathbf{W}'} [\langle \omega_{i1} - \omega'_{i1}, \tilde{\mathbf{x}} \rangle^2] + \frac{2}{m} \sum_{i=1}^m \mathbb{E}_{\mathbf{W}, \mathbf{W}'} [\langle \omega_{i2} - \omega'_{i2}, \tilde{\mathbf{x}} \rangle^2]. \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{\mathbf{W}, \mathbf{W}'} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}}) - \hat{y}(\mathbf{x}; \mu_{\mathbf{W}'}))^2]] &\leq \frac{2 \|\tilde{\Sigma}\|}{m} \mathbb{E}_{\mathbf{W}, \mathbf{W}'} [\|\mathbf{W} - \mathbf{W}'\|_F^2] \\ &= \frac{2 \|\tilde{\Sigma}\|}{m} W_2^2(\mu_t^m, \mu_\beta^{*\otimes m}). \end{aligned}$$

For the second term, notice that $\hat{y}(\mathbf{x}; \mu_\beta^*) = \mathbb{E}_{\mathbf{W}'} [\hat{y}(\mathbf{x}; \mu_{\mathbf{W}'})] = \mathbb{E}_{\mathbf{w}'_i} [\Psi(\mathbf{x}; \mathbf{w}'_i)]$ for all $1 \leq i \leq m$. By independence of (\mathbf{w}'_i) and Jensen's inequality, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{W}'} [(\hat{y}(\mathbf{x}; \mu_{\mathbf{W}'}) - \hat{y}(\mathbf{x}; \mu_\beta^*))^2] &= \frac{1}{m} \mathbb{E}_{\mathbf{w}'} [(\Psi(\mathbf{x}; \mathbf{w}') - \hat{y}(\mathbf{x}; \mu_\beta^*))^2] \\ &= \frac{1}{m} \mathbb{E}_{\mathbf{w}'} \left[\left(\int (\Psi(\mathbf{x}; \mathbf{w}') - \Psi(\mathbf{x}; \mathbf{w})) d\mu_\beta^*(\mathbf{w}) \right)^2 \right] \\ &\lesssim \frac{t^2}{m}. \end{aligned}$$

■

Thus, the rest of this section deals with establishing convergence rates for $\mathcal{F}_{\beta, \lambda}^m(\mu_l^m) \rightarrow \mathcal{F}_{\beta, \lambda}(\mu_\beta^*)$. To use the one-step decay of optimality gap provided by [44], we depend on the following assumption.

Assumption 3 *Suppose there exist constants L , C_L , and R , such that*

1. **(Lipschitz gradients of the Gibbs potential)** *For all $\mu, \mu' \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ and $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^{2d+2}$,*

$$\left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) - \nabla \hat{\mathcal{J}}'_0[\mu'](\mathbf{w}') \right\| \leq L(W_2(\mu, \mu') + \|\mathbf{w} - \mathbf{w}'\|), \quad (\text{B.10})$$

where W_2 is the 2-Wasserstein distance.

2. **(Bounded gradients of the Gibbs potential)** *For all $\mu \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ and $\mathbf{w} \in \mathbb{R}^{2d+2}$, we have $\left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) \right\| \leq R$.*
3. **(Bounded second variation)** *Denote the second variation of $\hat{\mathcal{J}}_0(\mu)$ at \mathbf{w} via $\hat{\mathcal{J}}''_0[\mu](\mathbf{w}, \mathbf{w}')$, which is defined as the first variation of $\mu \mapsto \hat{\mathcal{J}}'_0[\mu](\mathbf{w})$. Then, for all $\mu \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ and $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^{2d+2}$,*

$$\left| \hat{\mathcal{J}}''_0[\mu](\mathbf{w}, \mathbf{w}') \right| \leq L(1 + C_L(\|\mathbf{w}\|^2 + \|\mathbf{w}'\|^2)). \quad (\text{B.11})$$

We can now state the one-step bound.

Theorem 21 ([44, Theorem 2]) *Suppose $\hat{\mathcal{J}}_0$ satisfies Assumption 3. Assume $\lambda \lesssim 1$, $\beta, L, R \gtrsim 1$, and the initialization satisfies $\mathbb{E}[\|\mathbf{w}_0^i\|^2] \lesssim R^2$ for all $1 \leq i \leq m$. Then, for all $\eta \leq 1/4$,*

$$\mathcal{F}_{\beta, \lambda}^m(\mu_{l+1}^m) - \mathcal{F}_{\beta, \lambda}(\mu_\beta^*) \leq \exp\left(\frac{-\eta}{2\beta C_{\text{LSI}}}\right) (\mathcal{F}_{\beta, \lambda}^m(\mu_l^m) - \mathcal{F}_{\beta, \lambda}(\mu_\beta^*)) + \eta A_{m, \beta, \lambda, \eta}, \quad (\text{B.12})$$

where

$$A_{m,\beta,\lambda,\eta} := C \left(L^2 \left(d + \frac{R^2}{\lambda} \right) (\eta^2 + \frac{\eta}{\beta}) + \frac{L}{m\beta} \left(\frac{1}{C_{\text{LSI}}} + \left(\frac{R^2}{\lambda^2} + \frac{d}{\lambda\beta} \right) \left(\frac{C_L}{C_{\text{LSI}}} + \frac{L}{\beta} \right) \right) \right) \quad (\text{B.13})$$

for some absolute constant $C > 0$.

We now focus on bounding the constants that appear in Assumption 3.

Lemma 22 (Lipschitzness of $\nabla \hat{\mathcal{J}}'_0$) *Suppose ρ is either the squared error loss or is C_ρ Lipschitz and has a C'_ρ Lipschitz derivative. Assume $\kappa \gtrsim 1$. Notice that for the squared error loss, $C'_\rho = 1$. Then, for all $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ and $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^{2d+2}$, we have*

$$\left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) - \hat{\mathcal{J}}'_n[\mu'](\mathbf{w}') \right\| \lesssim \kappa C_\rho \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\| \|\mathbf{w} - \mathbf{w}'\| + C'_\rho \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\| W_2(\mu, \mu'),$$

for the Lipschitz loss, and

$$\left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) - \hat{\mathcal{J}}'_n[\mu'](\mathbf{w}') \right\| \lesssim \kappa \sqrt{\hat{\mathcal{J}}_0(\mu) \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}}^{\otimes 4}] \right\|_{2 \rightarrow 2}} \|\mathbf{w} - \mathbf{w}'\| + \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\| W_2(\mu, \mu'),$$

for the squared error loss, where $\left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}}^{\otimes 4}] \right\|_{2 \rightarrow 2} := \sup_{\|\mathbf{v}\| \leq 1} \left\| \mathbb{E}_{S_n} [\langle \tilde{\mathbf{x}}, \mathbf{v} \rangle^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|$.

Proof Recall that $\hat{\mathcal{J}}'_0[\mu](\mathbf{w}) = \mathbb{E}_{S_n} [\rho'(\hat{y}(\mathbf{x}; \mu) - y) \Psi(\mathbf{x}; \mathbf{w})]$, where $\Psi(\mathbf{x}; \mathbf{w}) = \phi_{\kappa,\iota}(\langle \boldsymbol{\omega}_1, \tilde{\mathbf{x}} \rangle) - \phi_{\kappa,\iota}(\langle \boldsymbol{\omega}_2, \tilde{\mathbf{x}} \rangle)$. We start with the triangle inequality,

$$\left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) - \nabla \hat{\mathcal{J}}'_0[\mu'](\mathbf{w}') \right\| \leq \left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) - \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}') \right\| + \left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}') - \nabla \hat{\mathcal{J}}'_0[\mu'](\mathbf{w}') \right\|.$$

We now focus on the first term. For the Lipschitz loss,

$$\begin{aligned} \left\| \nabla_{\boldsymbol{\omega}_1} \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) - \nabla_{\boldsymbol{\omega}_1} \hat{\mathcal{J}}'_0[\mu](\mathbf{w}') \right\| &= \left\| \mathbb{E}_{S_n} [\rho'(\hat{y}(\mathbf{x}; \mu) - y) (\phi'_\kappa(\langle \boldsymbol{\omega}_1, \tilde{\mathbf{x}} \rangle) - \phi'_\kappa(\langle \boldsymbol{\omega}'_1, \tilde{\mathbf{x}} \rangle)) \tilde{\mathbf{x}}] \right\| \\ &\leq C_\rho \mathbb{E}_{S_n} [(\phi'_\kappa(\langle \boldsymbol{\omega}_1, \tilde{\mathbf{x}} \rangle) - \phi'_\kappa(\langle \boldsymbol{\omega}'_1, \tilde{\mathbf{x}} \rangle))^2]^{1/2} \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2} \\ &\leq C_\rho \kappa \mathbb{E}_{S_n} [\langle \boldsymbol{\omega}_1 - \boldsymbol{\omega}'_1, \tilde{\mathbf{x}} \rangle^2]^{1/2} \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2} \\ &\leq C_\rho \kappa \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\| \|\boldsymbol{\omega}_1 - \boldsymbol{\omega}'_1\|, \end{aligned}$$

where the first inequality follows from Lemma 12, and the second inequality follows from the fact that $|\phi''_\kappa| \leq \kappa$. For the squared error loss, we have

$$\begin{aligned} \left\| \nabla_{\boldsymbol{\omega}_1} \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) - \nabla_{\boldsymbol{\omega}_1} \hat{\mathcal{J}}'_0[\mu](\mathbf{w}') \right\| &= \left\| \mathbb{E}_{S_n} [(\hat{y}(\mathbf{x}; \mu) - y) (\phi'_\kappa(\langle \mathbf{w}, \tilde{\mathbf{x}} \rangle) - \phi'_\kappa(\langle \mathbf{w}', \tilde{\mathbf{x}} \rangle)) \tilde{\mathbf{x}}] \right\| \\ &= \sup_{\|\mathbf{v}\| \leq 1} \mathbb{E}_{S_n} [(\hat{y}(\mathbf{x}; \mu) - y) (\phi'_\kappa(\langle \boldsymbol{\omega}_1, \tilde{\mathbf{x}} \rangle) - \phi'_\kappa(\langle \boldsymbol{\omega}'_1, \tilde{\mathbf{x}} \rangle)) \langle \mathbf{v}, \tilde{\mathbf{x}} \rangle] \\ &\leq \sup_{\|\mathbf{v}\| \leq 1} \sqrt{\mathbb{E}_{S_n} [(\hat{y}(\mathbf{x}; \mu) - y)^2] \mathbb{E}_{S_n} [(\phi'_\kappa(\langle \boldsymbol{\omega}_1, \tilde{\mathbf{x}} \rangle) - \phi'_\kappa(\langle \boldsymbol{\omega}'_1, \tilde{\mathbf{x}} \rangle))^2 \langle \mathbf{v}, \tilde{\mathbf{x}} \rangle^2]} \\ &\leq \kappa \sqrt{\hat{\mathcal{J}}_0(\mu) \sup_{\|\mathbf{v}\| \leq 1} \langle \mathbf{v}, \mathbb{E}_{S_n} [\langle \boldsymbol{\omega}_1 - \boldsymbol{\omega}'_1, \tilde{\mathbf{x}} \rangle^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \mathbf{v} \rangle} \\ &\leq \kappa \sqrt{\hat{\mathcal{J}}_0(\mu) \left\| \mathbb{E}_{S_n} [\langle \boldsymbol{\omega}_1 - \boldsymbol{\omega}'_1, \tilde{\mathbf{x}} \rangle^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|} \\ &\leq \kappa \sqrt{\hat{\mathcal{J}}_0(\mu) \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}}^{\otimes 4}] \right\|_{2 \rightarrow 2}} \|\boldsymbol{\omega}_1 - \boldsymbol{\omega}'_1\|. \end{aligned}$$

Similar bounds apply to the gradient with respect to ω_2 , which completes the bound on the first term of the triangle inequality.

We now consider the second term of the triangle inequality. Here we consider Lipschitz losses and the squared error loss at the same time since both have a Lipschitz derivative.

$$\begin{aligned}
 \left\| \nabla_{\omega_1} \hat{\mathcal{J}}'_0[\mu](\omega') - \nabla_{\omega_1} \hat{\mathcal{J}}'_0[\mu](\omega) \right\| &= \left\| (\rho'(\hat{y}(\mathbf{x}; \mu) - y) - \rho'(\hat{y}(\mathbf{x}; \mu') - y)) \phi'_\kappa(\langle \omega'_1, \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}} \right\| \\
 &\leq \mathbb{E}_{S_n} \left[(\rho'(\hat{y}(\mathbf{x}; \mu) - y) - \rho'(\hat{y}(\mathbf{x}; \mu') - y))^2 \right]^{1/2} \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2} \\
 &\leq C'_\rho \mathbb{E}_{S_n} [(\hat{y}(\mathbf{x}; \mu) - \hat{y}(\mathbf{x}; \mu'))^2]^{1/2} \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2},
 \end{aligned} \tag{B.14}$$

where the first inequality follows from Lemma 12. Let $\gamma \in \mathcal{P}_2(\mathbb{R}^{2d+2} \times \mathbb{R}^{2d+2})$ be a coupling of μ and μ' (i.e. the first and second marginals of γ are equal to μ and μ' respectively). Recall that,

$$\hat{y}(\mathbf{x}; \mu) - \hat{y}(\mathbf{x}; \mu') = \int (\phi_{\kappa, \iota}(\langle \omega_1, \tilde{\mathbf{x}} \rangle) - \phi_{\kappa, \iota}(\langle \omega_2, \tilde{\mathbf{x}} \rangle) - \phi_{\kappa, \iota}(\langle \omega'_1, \tilde{\mathbf{x}} \rangle) + \phi_{\kappa, \iota}(\langle \omega'_2, \tilde{\mathbf{x}} \rangle)) d\gamma(\mathbf{w}, \mathbf{w}').$$

Therefore by the triangle inequality for the L_2 norm $\mathbb{E}_{S_n} [(\cdot)^2]^{1/2}$ and Jensen's inequality,

$$\begin{aligned}
 \mathbb{E}_{S_n} [(\hat{y}(\mathbf{x}; \mu) - \hat{y}(\mathbf{x}; \mu'))^2]^{1/2} &\leq \mathbb{E}_{S_n} \left[\int (\phi_{\kappa, \iota}(\langle \omega_1, \tilde{\mathbf{x}} \rangle) - \phi_{\kappa, \iota}(\langle \omega'_1, \tilde{\mathbf{x}} \rangle))^2 d\gamma \right]^{1/2} \\
 &\quad + \mathbb{E}_{S_n} \left[\int (\phi_{\kappa, \iota}(\langle \omega_2, \tilde{\mathbf{x}} \rangle) - \phi_{\kappa, \iota}(\langle \omega'_2, \tilde{\mathbf{x}} \rangle))^2 d\gamma \right]^{1/2} \\
 &\leq \int \mathbb{E}_{S_n} [\langle \omega_1 - \omega'_1, \tilde{\mathbf{x}} \rangle^2]^{1/2} d\gamma + \int \mathbb{E}_{S_n} [\langle \omega_2 - \omega'_2, \tilde{\mathbf{x}} \rangle^2]^{1/2} d\gamma \\
 &\leq \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2} \int (\|\omega_1 - \omega'_1\| + \|\omega_2 - \omega'_2\|) d\gamma(\mathbf{w}_1, \mathbf{w}_2) \\
 &\leq \sqrt{2 \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|} \int \|\mathbf{w} - \mathbf{w}'\|^2 d\gamma(\mathbf{w}, \mathbf{w}').
 \end{aligned}$$

By choosing γ whose transport cost attains (or converges to) the optimal cost, we have

$$\mathbb{E}_{S_n} [(\hat{y}(\mathbf{x}; \mu) - \hat{y}(\mathbf{x}; \mu'))^2]^{1/2} \leq \sqrt{2 \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|} W_2(\mu, \mu').$$

Plugging the above result into (B.14), we have

$$\left\| \nabla_{\omega_1} \hat{\mathcal{J}}'_0[\mu](\omega') - \nabla_{\omega_1} \hat{\mathcal{J}}'_0[\mu](\omega) \right\| \leq \sqrt{2} C'_\rho \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2} W_2(\mu, \mu').$$

Notice that the same bound holds for gradients with respect to ω_2 . Thus the bound of the second term in the triangle inequality and the proof is complete. \blacksquare

Lemma 23 (Boundedness of $\nabla \hat{\mathcal{J}}'_0$) *In the same setting as Lemma 22, for all $\mu \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ and $\mathbf{w} \in \mathbb{R}^{2d+2}$, we have*

$$\left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) \right\| \leq \sqrt{2} \tilde{C}_\rho \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2},$$

where $\tilde{C}_\rho = C_\rho$ when ρ is Lipschitz and $\tilde{C}_\rho = \sqrt{2 \hat{\mathcal{J}}_0(\mu)}$ when ρ is the squared error loss.

Proof Notice that $|\phi'_\kappa| \leq 1$. Therefore,

$$\begin{aligned} \left\| \nabla_{\omega_1} \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) \right\| &= \left\| \mathbb{E}_{S_n} [\rho'(\hat{y} - y) \phi'_\kappa(\langle \omega_1, \tilde{\mathbf{x}} \rangle) \tilde{\mathbf{x}}] \right\| \\ &\leq \sqrt{\mathbb{E}_{S_n} [\rho'(\hat{y} - y)^2] \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|} \\ &\leq \tilde{C}_\ell \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2}, \end{aligned}$$

where the first inequality follows from Lemma 12. \blacksquare

Lemma 24 (Boundedness of $\hat{\mathcal{J}}''_0$) *In the same setting as Lemma 22, for all $\mu \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ and $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^{2d+2}$, we have*

$$\left| \hat{\mathcal{J}}''_0[\mu](\mathbf{w}, \mathbf{w}') \right| \leq C'_\rho \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\| \left(\|\mathbf{w}\|^2 + \|\mathbf{w}'\|^2 \right),$$

where we recall that $C'_\rho = 1$ for the squared error loss.

Proof It is straightforward to show that

$$\hat{\mathcal{J}}''_0[\mu](\mathbf{w}, \mathbf{w}') = \mathbb{E}_{S_n} [\rho''(\hat{y}(\mathbf{x}; \mu) - y) \Psi(\mathbf{x}; \mathbf{w}) \Psi(\mathbf{x}; \mathbf{w}')].$$

Then, by the Cauchy-Schwartz inequality,

$$\hat{\mathcal{J}}''_0[\mu](\mathbf{w}, \mathbf{w}') \leq C'_\rho \mathbb{E}_{S_n} [\Psi(\mathbf{x}; \mathbf{w})^2]^{1/2} \mathbb{E}_{S_n} [\Psi(\mathbf{x}; \mathbf{w}')^2]^{1/2}.$$

Moreover, by the Lipschitzness of $\phi_{\kappa, \iota}$,

$$\begin{aligned} \mathbb{E}_{S_n} [\Psi(\mathbf{x}; \mathbf{w})^2]^{1/2} &\leq \mathbb{E}_{S_n} [\langle \omega_1, \tilde{\mathbf{x}} \rangle^2]^{1/2} + \mathbb{E}_{S_n} [\langle \omega_2, \tilde{\mathbf{x}} \rangle^2]^{1/2} \\ &\leq \sqrt{2} \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\|^{1/2} \|\mathbf{w}\| \end{aligned}$$

We can similarly bound the expression for \mathbf{w}' , and arrive at the statement of the lemma via Young's inequality,

$$\hat{\mathcal{J}}''_0[\mu](\mathbf{w}, \mathbf{w}') \leq 2 \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\| C'_\rho \|\mathbf{w}\| \|\mathbf{w}'\| \leq \left\| \mathbb{E}_{S_n} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top] \right\| (\|\mathbf{w}\|^2 + \|\mathbf{w}'\|^2). \quad \blacksquare$$

In particular, the gradient bound of Lemma 23 implies the following LSI estimate, which follows from [44, Theorem 1].

Proposition 25 ([44, Theorem 1]) *Suppose $\left\| \nabla \hat{\mathcal{J}}'_0[\mu](\mathbf{w}) \right\| \leq R$ for all $\mu \in \mathcal{P}_2$ and $\mathbf{w} \in \mathbb{R}^{2d+2}$. Then, the family of probability measures $\nu_\mu \propto \exp(-\beta G'[\mu])$ for $\mu \in \mathcal{P}_2(\mathbb{R}^{2d+2})$ satisfy a uniform LSI with constant*

$$C_{\text{LSI}} \lesssim \frac{1}{\beta \lambda} \exp\left(\frac{4\beta R^2 \sqrt{2d/\pi}}{\lambda}\right) \wedge \left\{ \frac{1}{\beta \lambda} + \exp\left(\frac{\beta R^2}{2\lambda}\right) \left(\frac{R^2}{\lambda^2} + \frac{1}{\beta \lambda}\right) \left(d + \frac{\beta R^2}{\lambda}\right) \right\}. \quad (\text{B.15})$$

We collect the smoothness estimates and simplify them under the event of Lemma 10 in the following Corollary.

Corollary 26 *Suppose ρ and ρ' are C_ρ and C'_ρ Lipschitz respectively, with $C_\rho, C'_\rho \lesssim 1$. Recall that $\Sigma := \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$. On the event of Lemma 10, we have $\|\mathbb{E}_{S_n}[\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top]\| \lesssim \|\Sigma\| \vee \tilde{r}_x^2$, and consequently, $\hat{\mathcal{J}}_0^l$ satisfies Assumption 3 with constants $L \lesssim \kappa(\|\Sigma\| \vee \tilde{r}_x^2)$, $R \lesssim \|\Sigma\|^{1/2} \vee \tilde{r}_x$, and $C_L = \kappa^{-1}$.*

Using the estimates above, we can present the following convergence bound $\mathcal{F}_{\beta,\lambda}^m(\mu_\beta^*) - \mathcal{F}_{\beta,\lambda}(\mu_\beta^*)$.

Proposition 27 *Let $\tilde{r}_x := \|\Sigma\| \vee \tilde{r}_x$, and for simplicity assume $C_{\text{LSI}} \geq \beta$. For any $\varepsilon \lesssim 1$, suppose the step size satisfies*

$$\eta \lesssim \frac{\varepsilon}{C_{\text{LSI}}\kappa^2\tilde{r}_x^4(d + \tilde{r}_x^2/\lambda)},$$

the width of the network satisfies,

$$m \gtrsim \frac{\kappa\tilde{r}_x^2 \left(1 + \left(\frac{\tilde{r}_x^2}{\lambda^2} + \frac{d}{\lambda\beta}\right) \left(\frac{1}{\kappa} + \frac{\kappa\tilde{r}_x^2 C_{\text{LSI}}}{\beta}\right)\right)}{\varepsilon},$$

and the number of iterations satisfies

$$l \gtrsim \frac{\beta C_{\text{LSI}}}{\eta} \ln \left(\frac{\mathcal{F}_{\beta,\lambda}^m(\mu_0^m) - \mathcal{F}_{\beta,\lambda}^*}{\varepsilon} \right).$$

Then, we have $\mathcal{F}_{\beta,\lambda}^m(\mu_l^m) - \mathcal{F}_{\beta,\lambda}(\mu_\beta^*) \leq \varepsilon$.

Proof Throughout the proof, we will assume the event of Lemma 10 holds. Let $\mathcal{F}_{\beta,\lambda}^* := \mathcal{F}_{\beta,\lambda}(\mu_\beta^*)$. Notice that by iterating the bound of Theorem 21, we have

$$\begin{aligned} \mathcal{F}_{\beta,\lambda}^m(\mu_l^m) - \mathcal{F}_{\beta,\lambda}^* &\leq \exp\left(\frac{-l\eta}{2\beta C_{\text{LSI}}}\right) (\mathcal{F}_{\beta,\lambda}^m(\mu_0^m) - \mathcal{F}_{\beta,\lambda}^*) + \frac{\eta A_{m,\beta,\lambda,\eta}}{1 - \exp\left(\frac{-\eta}{2\beta C_{\text{LSI}}}\right)} \\ &\leq \exp\left(\frac{-l\eta}{2\beta C_{\text{LSI}}}\right) (\mathcal{F}_{\beta,\lambda}^m(\mu_0^m) - \mathcal{F}_{\beta,\lambda}^*) + 4\beta C_{\text{LSI}} A_{m,\beta,\lambda,\eta}, \end{aligned}$$

where the second inequality holds for $\eta \leq 2\beta C_{\text{LSI}}$ since $1 - e^{-x} \geq x/2$ for $x \in [0, 1]$. We now bound $A_{m,\beta,\lambda,\eta}$ so that the RHS of the above is less than $\mathcal{O}(\varepsilon)$ by choosing a sufficiently large m and a sufficiently small η . Recall that given constants L and R from Assumption 3,

$$A_{m,\beta,\lambda,\eta} \asymp L^2 \left(d + \frac{R^2}{\lambda}\right) \left(\eta^2 + \frac{\eta}{\beta}\right) + \frac{L}{m\beta} \left(\frac{1}{C_{\text{LSI}}} + \left(\frac{R^2}{\lambda^2} + \frac{d}{\lambda\beta}\right) \left(\frac{C_L}{C_{\text{LSI}}} + \frac{L}{\beta}\right)\right).$$

From Corollary 26, $L \asymp \kappa(\|\Sigma\| \vee \tilde{r}_x^2)$, $R \asymp \|\Sigma\|^{1/2} \vee \tilde{r}_x$, and $C_L = \kappa^{-1}$. To avoid notational clutter, let $\tilde{r}_x^2 := \|\Sigma\| \vee \tilde{r}_x^2$. Then, to control the terms containing η , it suffices to choose

$$\eta \lesssim \sqrt{\frac{\varepsilon}{\beta C_{\text{LSI}}\kappa^2\tilde{r}_x^4(d + \tilde{r}_x^2/\lambda)}} \wedge \frac{\varepsilon}{C_{\text{LSI}}\kappa^2\tilde{r}_x^4(d + \tilde{r}_x^2/\lambda)},$$

for which we can simply choose

$$\eta \lesssim \frac{\varepsilon}{C_{\text{LSI}}\kappa^2\tilde{r}_x^4(d + \tilde{r}_x^2/\lambda)}.$$

Further, to control the term containing the number of particles m , we need

$$m \gtrsim \frac{\kappa \tilde{r}_x^2 \left(1 + \left(\frac{\tilde{r}_x^2}{\lambda^2} + \frac{d}{\lambda \beta} \right) \left(\frac{1}{\kappa} + \frac{\kappa \tilde{r}_x^2 C_{\text{LSI}}}{\beta} \right) \right)}{\varepsilon}.$$

To drive the suboptimality bound below ε , we also need to let the number of iterations l satisfy

$$l \gtrsim \frac{\beta C_{\text{LSI}}}{\eta} \ln \left(\frac{\mathcal{F}_{\beta, \lambda}^m(\mu_0^m) - \mathcal{F}_{\beta, \lambda}^*}{\varepsilon} \right).$$

With the above conditions, we can guarantee

$$\mathcal{F}_{\beta, \lambda}^m(\mu_l^m) - \mathcal{F}_{\beta, \lambda}^* \lesssim \varepsilon,$$

which finishes the proof. \blacksquare

Finally, we are ready to present the proof of Theorem. Specifically, we will prove the following theorem which is the more detailed version of Theorem 2. To do so, we introduce the following assumption, which can be verified by Proposition 25.

Assumption 4 Let $\mathbf{W}^l = (\mathbf{w}_1^l, \dots, \mathbf{w}_m^l)$ denote the trajectory of MFLA. Then, the probability measure $\nu_{\mu_{\mathbf{W}^l}} \propto \exp(-\beta \hat{\mathcal{J}}'_\lambda[\mu_{\mathbf{W}^l}])$ satisfies the LSI (A.6) with constant C_{LSI} for all $l \geq 0$. For simplicity, assume $C_{\text{LSI}} \geq \beta$.

With the above assumption, we can state the detailed version of Theorem 2.

Theorem 28 Under Assumptions 1, 2 and 4, consider MFLA (2.3) with parameters $\beta = \tilde{\Theta}(d_{\text{eff}})$, $\lambda = \tilde{\lambda} r_x^2$ and $\eta \leq \tilde{\mathcal{O}}\left(\frac{1}{C_{\text{LSI}} \kappa^2 d}\right)$. The algorithm is initialized with the weights sampled i.i.d. from some distribution $\mathbf{w}_j^0 \sim \mu_0$. Suppose $\tilde{\lambda}, \kappa^{-1} = o_n(1)$, and ι is chosen via the lower bound of Lemma 19. Then, with the number of samples n , the number of neurons m , and the number of iterations l that can respectively be bounded by

$$n \leq \tilde{\mathcal{O}}(d_{\text{eff}}), \quad m \leq \tilde{\mathcal{O}}\left(\frac{C_{\text{LSI}} \kappa^2 d}{\beta^2 \lambda}\right), \quad l \leq \tilde{\mathcal{O}}\left(\frac{C_{\text{LSI}} \beta}{\eta}\right), \quad (\text{B.16})$$

with probability at least $1 - \mathcal{O}(n^{-q})$ for some $q > 0$, the excess risk satisfies

$$\mathbb{E}_{\mathbf{W}^l} \mathbb{E}_{y, \mathbf{x}}[\rho(y - \hat{y}_m(\mathbf{x}; \mathbf{W}^l))] - \mathbb{E}_\xi[\rho(\xi)] \leq o_n(1). \quad (\text{B.17})$$

For the statement of Theorem 2, we can choose $\tilde{\lambda}^{-1}, \kappa = \text{polylog}(n)$.

B.5. Proof of Theorem 2

Recall that $\lambda = \tilde{\lambda} r_x^2$, and let $\beta = \frac{d_{\text{eff}} + \tilde{r}_x^2 / r_x^2}{\varepsilon^2 \tilde{\lambda}}$ and $n \geq \frac{(d_{\text{eff}} + \tilde{r}_x^2 / r_x^2) \iota^2}{\tilde{\lambda} \varepsilon^4}$ for some $\varepsilon \lesssim 1$, where $\tilde{\varepsilon} := \tilde{\mathcal{O}}(\tilde{\lambda}^{\frac{1}{k+2}} + \varepsilon + \kappa^{-1})$. Then, from Lemma 19, we have $\mathcal{J}_0(\mu_\beta^*) - \mathbb{E}[\rho(\xi)] \lesssim \tilde{\varepsilon}$. On the other hand, given the step size η , width m , and number of iterations l by Proposition 27, we have $\mathcal{F}_{\beta, \lambda}^m(\mu_l^m) - \mathcal{F}_{\beta, \lambda}(\mu_\beta^*) \leq \varepsilon$. Therefore,

$$\mathbb{E}_{\mathbf{W} \sim \mu_l^m} [J_0(\mathbf{W})] - \mathcal{J}_0(\mu_\beta^*) \lesssim \sqrt{\frac{\tilde{r}_x^2 \beta C_{\text{LSI}} \varepsilon}{m}} + \frac{\iota^2}{m}.$$

Additionally, from Lemma 19, we have

$$\beta^{-1} \mathcal{H}(\mu_\beta^* | \gamma) \lesssim \mathbb{E}[\rho(\xi)] + \tilde{\mathcal{O}}(\tilde{\lambda}^{\frac{1}{k+2}}) + \varepsilon + \kappa^{-1} \lesssim 1.$$

Consequently, for $m \geq \frac{\bar{r}_x^2 (d_{\text{eff}} + \bar{r}_x^2 / r_x^2) C_{\text{LSI}}}{\tilde{\lambda} \varepsilon^3} \vee \frac{L}{\varepsilon^2}$, we have $\mathbb{E}_{\mathbf{W} \sim \mu_l^m} [J_0(\mathbf{W})] - \mathcal{J}_0(\mu_\beta^*) \leq \varepsilon$. Therefore, combining the bounds above, we have

$$\mathbb{E}_{\mathbf{W} \sim \mu_l^m} [J_0(\mathbf{W})] - \mathbb{E}[\rho(\xi)] \lesssim \tilde{\mathcal{O}}(\tilde{\lambda}^{\frac{1}{k+2}}) + \varepsilon + \kappa^{-1}.$$

Consequently, we can take $\tilde{\lambda} = o_n(1)$, $\varepsilon = o_n(1)$, $\kappa^{-1} = o_n(1)$ to finish the proof.

Appendix C. Comparisons with the Formulation of [39]

Here, we provide a number of comparisons with results of [39]. In Section C.1, we show that the statistical model (2.1) is more general than their formulation, even for parity learning problems. In Section C.2, we provide an informal comparison of their effective dimension to our setting, exhibiting the improvement in our definition of effective dimension.

C.1. Generality of the Formulations

We begin by pointing out that the formulation of k -index model of (2.1) is strictly more general than that of [39], even for learning k -sparse parities. Recall that in their setting, they consider inputs of the type $\mathbf{x} = \Sigma^{1/2} \mathbf{z}$ for some positive definite Σ , where $\mathbf{z} \sim \text{Unif}(\{\pm 1\}^d)$ (their original formulation uses $\mathbf{z} \sim \text{Unif}(\{\pm 1/\sqrt{d}\}^d)$, but we rescale the input to be consistent with the notation of this paper). The labels are given by

$$y = \text{sign} \left(\prod_{i=1}^k \langle \tilde{\mathbf{u}}_i, \mathbf{z} \rangle \right) = \text{sign} \left(\prod_{i=1}^k \langle \Sigma^{-1/2} \tilde{\mathbf{u}}_i, \mathbf{x} \rangle \right), \quad (\text{C.1})$$

where $\{\tilde{\mathbf{u}}_i\}_{i=1}^k$ are orthonormal vectors. Then, we can define an orthonormal set of vectors $\{\mathbf{u}_i\}_{i=1}^k$ such that $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k) = \text{span}(\Sigma^{-1/2} \tilde{\mathbf{u}}_1, \dots, \Sigma^{-1/2} \tilde{\mathbf{u}}_k)$, and define g such that

$$g \left(\frac{\langle \mathbf{u}_1, \mathbf{x} \rangle}{\sqrt{k}}, \dots, \frac{\langle \mathbf{u}_k, \mathbf{x} \rangle}{\sqrt{k}} \right) = g \left(\frac{\langle \Sigma^{1/2} \mathbf{u}_1, \mathbf{z} \rangle}{\sqrt{k}}, \dots, \frac{\langle \Sigma^{1/2} \mathbf{u}_k, \mathbf{z} \rangle}{\sqrt{k}} \right) = \text{sign} \left(\prod_{i=1}^k \langle \tilde{\mathbf{u}}_i, \mathbf{z} \rangle \right),$$

for all $\mathbf{z} \in \{\pm 1\}^d$. Therefore, the parity formulation of (C.1) can be seen as a special case of the k -index model (2.1). Note that g is only defined on 2^d points, and we can extend it to all of \mathbb{R}^k such that $g : \mathbb{R}^k \rightarrow \mathbb{R}$ is Lipschitz.

In contrast, the k -index model can represent parity problems that cannot be represented by (C.1). Starting from an orthonormal set of vectors $\{\mathbf{u}_i\}_{i=1}^k$ in \mathbb{R}^d , let

$$y = g \left(\frac{\langle \mathbf{u}_1, \mathbf{x} \rangle}{\sqrt{k}}, \dots, \frac{\langle \mathbf{u}_k, \mathbf{x} \rangle}{\sqrt{k}} \right) = \text{sign} \left(\prod_{i=1}^k \langle \mathbf{u}_i, \mathbf{x} \rangle \right). \quad (\text{C.2})$$

Consider the case where $k = 2$, then $y = \text{sign} \left(\left\langle \Sigma^{1/2} \mathbf{u}_1, \mathbf{z} \right\rangle \left\langle \Sigma^{1/2} \mathbf{u}_2, \mathbf{z} \right\rangle \right)$. To be able to reformulate this to (C.1), we need to find orthonormal $\tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2 \in \mathbb{R}^d$ such that

$$\text{sign} \left(\left\langle \Sigma^{1/2} \mathbf{u}_1, \mathbf{z} \right\rangle \left\langle \Sigma^{1/2} \mathbf{u}_2, \mathbf{z} \right\rangle \right) = \text{sign}(\langle \tilde{\mathbf{u}}_1, \mathbf{z} \rangle \langle \tilde{\mathbf{u}}_2, \mathbf{z} \rangle), \quad \forall \mathbf{z} \in \{\pm 1\}^d.$$

If Σ has rank less than d such that $\Sigma^{1/2} \mathbf{u}_1 = \Sigma^{1/2} \mathbf{u}_2$, then the above implies $\text{sign}(\langle \tilde{\mathbf{u}}_1, \mathbf{z} \rangle \langle \tilde{\mathbf{u}}_2, \mathbf{z} \rangle) \geq 0$ for all $\mathbf{z} \in \{\pm 1\}^d$. In particular, we must have some \mathbf{z} where $\text{sign}(\langle \tilde{\mathbf{u}}_1, \mathbf{z} \rangle \langle \tilde{\mathbf{u}}_2, \mathbf{z} \rangle) > 0$, which implies that

$$\sum_{i=1}^{2^d} \langle \tilde{\mathbf{u}}_1, \mathbf{z}_i \rangle \langle \tilde{\mathbf{u}}_2, \mathbf{z}_i \rangle = \left\langle \tilde{\mathbf{u}}_1, \sum_{i=1}^{2^d} \mathbf{z}_i \mathbf{z}_i^\top \tilde{\mathbf{u}}_2 \right\rangle = 2^d \langle \tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2 \rangle > 0,$$

which is in contradiction with $\langle \tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2 \rangle = 0$. Therefore, for such Σ , we cannot formulate (C.2) as a special case of (C.1). This argument is robust with respect to small perturbations of Σ which make it full-rank, implying that (C.2) is strictly more general than (C.1), even when only considering full-rank covariance matrices.

C.2. Comparison with the Effective Dimension of [39]

A close inspection of the proofs in [39] demonstrates that one can define their effective dimension in a scale invariant manner as $\tilde{d}_{\text{eff}} := \text{tr}(\Sigma) \left\| \sum_{i=1}^k \Sigma^{-1/2} \tilde{\mathbf{u}}_i \right\|^2$. From the previous section, we observed that to reduce their setting to ours, we need to choose a set $\{\mathbf{u}_i\}_{i=1}^k$ of normalized vectors that spans the set of vectors $\{\Sigma^{-1/2} \tilde{\mathbf{u}}_i\}_{i=1}^k$. In particular, we can choose $\mathbf{u}_i = \frac{\Sigma^{-1/2} \tilde{\mathbf{u}}_i}{\|\Sigma^{-1/2} \tilde{\mathbf{u}}_i\|}$, or equivalently write $\tilde{\mathbf{u}}_i = \frac{\Sigma^{1/2} \mathbf{u}_i}{\|\Sigma^{1/2} \mathbf{u}_i\|}$. While $\{\mathbf{u}_i\}_{i=1}^k$ are not orthogonal, our proofs do not rely on the orthogonality assumption and it is only made for simplicity. Hence, we have

$$\tilde{d}_{\text{eff}} = \text{tr}(\Sigma) \left\| \sum_{i=1}^k \frac{\mathbf{u}_i}{\|\Sigma^{1/2} \mathbf{u}_i\|} \right\|^2 \leq k \text{tr}(\Sigma) \sum_{i=1}^k \left\| \Sigma^{1/2} \mathbf{u}_i \right\|^{-2}.$$

Note that the above upper bound is sharp when $k = 1$, and is lower bounded by our definition of effective dimension stated in Definition 1.