

Package ‘geometry’ detected.This can cause a problem for jmlrbook

Training Dynamics of Softmax Self-Attention: Fast Global Convergence via Preconditioning

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2026

Abstract

We study the training dynamics of gradient descent in a softmax self-attention layer trained to perform linear regression and propose a first-order optimization algorithm which converges to the globally optimal self-attention parameters at a geometric rate. Our analysis proceeds in two steps. First, we show that in the infinite-data limit the regression problem solved by the self-attention layer is equivalent to a nonconvex matrix factorization problem. Second, we exploit this connection to design a novel “structure-aware” variant of gradient descent which efficiently optimizes the original finite-data regression objective. Our optimization algorithm features several innovations over vanilla gradient descent, including a data-dependent preconditioner and a scale-invariant regularizer which help avoid spurious stationary points, a renormalization step which ensures that the softmax parameters remain bounded, and a spectral initialization of parameters which lie near the manifold of global minima with high probability. We prove that the generalization error of the model trained by our algorithm decreases exponentially fast in the number of gradient descent iterations, up to an additional error term that decreases as $1/n$ in the size of the training set.

1. Introduction

The self-attention mechanism is a neural architecture originally proposed by [3] for machine translation. It was subsequently adopted by [12] to form the basis of the Transformer architecture, which underlies many recent advances in natural language processing [10] and computer vision [6]. It has proven to be remarkably versatile, and can additionally be trained to mimic various algorithms from statistics, optimization, and machine learning [7]. Despite its numerous empirical successes, our theoretical understanding of the self-attention mechanism remains poor. Many of the prior works on the theoretical behavior of self-attention are conditional; they prove that *if* the self-attention parameters could be optimized to their globally optimal values then the resulting model would exhibit strong performance on various downstream tasks, but they do not establish *when* such optimization is possible or *how* it should be performed, e.g., [4, 9]. Analysis of the training dynamics of gradient descent in a softmax self-attention layer is exceedingly challenging due to the pairwise nonlinear interactions appearing in the softmax function.

A recent line of theoretical works (e.g., [1, 2, 5, 15]) seek to understand the optimization dynamics of self-attention in the setting of random linear regression¹ which was empirically investigated by [7] and [14]. In this model, the covariates are drawn from a distribution, and the response variables are a noisy linear function of the covariates. The performance of a predictor is measured using the square loss. A natural question is whether a self-attention mechanism can be trained to accurately predict the label corresponding to a given covariate; even in this simple setting, this is not at all obvious due to the nonconvexity of the loss in the model parameters. While several of these papers derive various global convergence guarantees,

1. The term ‘random linear regression’ is a misnomer, since we are studying a setting where a nonlinear model is used to fit data which is generated by a planted linear model (i.e., nonlinear regression). We adopt this terminology to remain consistent with preexisting literature.

they suffer from two drawbacks. First, most of these works only study a simplified, linearized variant of self-attention instead of the original softmax attention mechanism of [3]. Second, all of these works only study the optimization dynamics in an asymptotic limit where either the learning algorithm has access to infinitely many samples, or has an unlimited budget of gradient iterations to converge to optimality; none of these works quantify how model performance depends on the number of samples or the compute budget. In this paper we address both of these challenges.

1.1. Main Contributions

We consider a setting where the number of self-attention parameters is fixed, and bound the generalization error as a function of the size of the training set n and the number of gradient descent iterations m . Our contributions can be summarized as follows.

1. In Section 3, we study the population loss, i.e., the asymptotic limit of the training loss as n approaches infinity. We show that this loss has a simple closed-form description and show that it is equivalent to a certain weighted matrix factorization loss. Using ideas from the matrix factorization literature, we propose a novel regularizer of the population loss and show that the regularized population loss has infinitely many global minima which together form a smooth connected manifold. While this regularized loss is globally nonconvex, we prove that it exhibits one-point strong convexity and one-point smoothness near the manifold of global minima in a certain geometry in which the inner product between two points is weighted by the covariance of the data distribution.
2. In Section 4 we leverage the geometric results of Section 3 to design a “structure-aware” gradient descent algorithm which finds a set of model parameters with near-optimal population loss. Our algorithm features several innovations over standard optimization algorithms such as SGD and Adam. First, our algorithm initializes the parameters of the self-attention layer using spectral information from the training data; we show that these parameters lie near the manifold of global optima with high probability. Second, our algorithm incorporates a data-dependent regularizer and a data-dependent preconditioner which mirror the form of the regularizer and preconditioner described in Theorem ???. Our algorithm can be viewed as evolving each of the self-attention parameters in the geometry most natural to that parameter. We next present our main result: a scaling law which describes how the generalization error decreases as a function of n and m . We decompose the excess risk of the model trained by our optimization algorithm into two pieces: a statistical error term, which measures how well the random, finite-data empirical loss approximates the population loss, and an optimization error term, which measures the distance between the parameters found by our model and the globally optimal parameters. We show that the statistical error decreases at a fast $1/n$ rate, up to logarithmic factors, and the optimization error decays exponentially in m . To the best of our knowledge, this is the first result which establishes fast (i.e., geometric rate) global convergence of a first-order method on a softmax self-attention training objective in any setting, and also the first result which establishes a sharp $1/n$ statistical rate in any setting.

In the course of proving our main results, we establish novel concentration inequalities which bound how often the softmax-weighted mean and softmax-weighted covariance deviate far from their respective means; these results may be of independent interest.

2. Model

We study regression using the square loss, where the covariates are d -dimensional and the response variables are p -dimensional. Specifically, we consider a setting where we are given n samples $\{x_i, y_i\}_{i=1}^n$, where each x_i is drawn independently from $\mathcal{N}(0, \Sigma)$ and each response y_i has the form $y_i = Mx_i + z_i$ for some fixed

weight matrix $M \in \mathbb{R}^{p \times d}$. The noise variables $\{z_i\}_{i=1}^n$ are drawn i.i.d. and independently from the covariates from $\mathcal{N}(0, \Omega)$. Our goal is to learn a prediction rule which, when given a fresh covariate $x \sim \mathcal{N}(0, \Sigma)$, generates a prediction \hat{y} which is close to Mx . We consider the family of regression functions consisting of single-layer single-head softmax self-attention functions; such functions are parameterized by $\theta = (A, B)$, where $A \in \mathbb{R}^{p \times d}$ and $B \in \mathbb{R}^{d \times d}$. We think of θ as the vertical concatenation A and B , so that $\theta \in \mathbb{R}^{(p+d) \times d}$. Given a fresh covariate $x \in \mathbb{R}^d$, such functions predict a corresponding label \hat{y} given by

$$\hat{y} = A \left(\frac{\sum_{j=1}^n \exp(x^\top B x_j) x_j}{\sum_{j=1}^n \exp(x^\top B x_j)} \right).$$

The prediction \hat{y} is the image of a convex combination of the covariates $\{x_j\}_{j=1}^n$ under the linear map A , where the weights of this convex combination are determined by the nonlinear softmax function parameterized by B . We note that in the language of the original self-attention paper [12], the parameter A is called the value matrix, while B is the product of the key and query matrices.

We define the in-sample empirical loss

$$\hat{L}(\theta) = \frac{1}{2n} \sum_{i=1}^n \left\| A \left(\frac{\sum_{j=1}^n \exp(x_i^\top B x_j) x_j}{\sum_{j=1}^n \exp(x_i^\top B x_j)} \right) - y_i \right\|_2^2. \quad (1)$$

We also define the population loss

$$L(\theta) = \frac{1}{2} \mathbb{E}_{x_1, z_1} \left\| A \left(\frac{\mathbb{E}_{x_2} [\exp(x_1^\top B x_2) x_2]}{\mathbb{E}_{x_2} [\exp(x_1^\top B x_2)]} \right) - (M x_1 + z_1) \right\|_2^2, \quad (2)$$

where x_1 and x_2 are sampled independently from $\mathcal{N}(0, \Sigma)$ and z_1 is sampled from $\mathcal{N}(0, \Omega)$. The population loss has the following intuitive interpretation. When n is large, we expect that each of the summations appearing in the numerator and the denominator of the predictor \hat{y} should approach their respective expectations. Averaging over the individual losses, we obtain the population loss.

3. Structure of the population loss

We characterize the population loss $L(\theta)$ in closed form, and show that a regularized variant $Q(\theta) = L(\theta) + R(\theta)$ of the population loss obeys certain convexity and smoothness properties near its minima. The regularizer $R(\theta)$ has the interpretation of ‘‘balancing’’ A and B , so as to make sure that the optimization algorithm does not spend too much effort optimizing one parameter at the expense of the other. A crucial aspect of our result is that these properties hold only in a certain geometry in which the inner product and norm are reweighted by the covariance Σ . As a consequence, our convexity and smoothness results are not stated in terms of the raw gradient $\nabla Q(\theta)$, but rather in terms of the preconditioned gradient $P^{-1} \nabla Q$, where P^{-1} is a preconditioner depending on Σ . A complete characterization of the population loss appears in Section A.2.

4. Main Result

We propose a first-order algorithm which converges to the population-optimal self-attention parameters at a geometric rate. This algorithm is formally described in the display Algorithm 1. We have shown in Section A.2 that the regularized loss $Q(\theta)$ is strongly convex and smooth near the manifold of global minima \mathcal{S} , provided that the gradient of $\nabla Q(\theta)$ is appropriately preconditioned. A natural idea is to initialize the parameters near \mathcal{S} using empirical estimates of $\hat{\Sigma}$ and \hat{M} , and then to approximate $P^{-1} \nabla Q(\theta)$ by $\hat{P}^{-1}(\hat{L}(\theta) + \hat{R}(\theta))$, where $\hat{R}(\theta)$ is simply the regularizer obtained by replacing the true covariance Σ by its

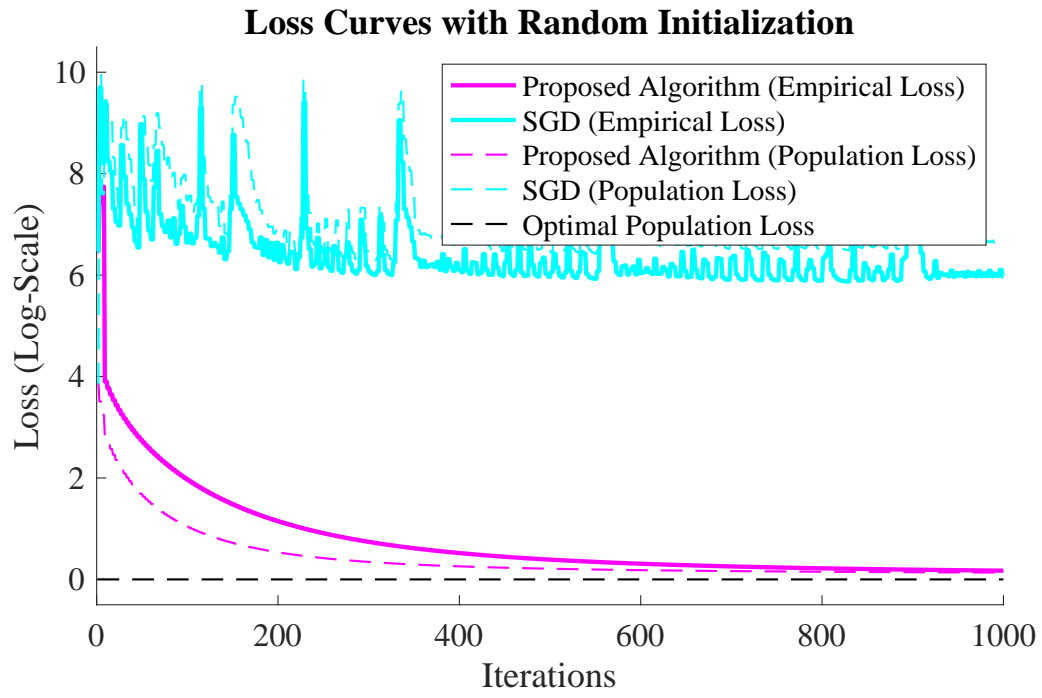


Figure 1: We plot the training losses and test losses incurred by our algorithm and SGD, where both algorithms are initialized at the same random point.

Algorithm 1 Preconditioned Gradient Descent for Self-Attention

Input: Data $\{(x_i, y_i)\}_{i=1}^n$, step size $\eta > 0$, iteration budget m , failure probability δ
Output: (A_m, B_m)

$$\hat{\Sigma} \leftarrow n^{-1} \sum_{i=1}^n x_i x_i^\top, \quad \hat{M} \leftarrow n^{-1} \sum_{i=1}^n y_i x_i^\top \hat{\Sigma}^{-1}$$

$$\hat{U} \hat{\Gamma} \hat{V}^\top \leftarrow \text{SVD}(\hat{M} \hat{\Sigma}^{1/2}), \quad A_0 \leftarrow \hat{U} \hat{\Gamma}^{1/2} \hat{\Sigma}^{-1/2}, \quad B_0 \leftarrow \hat{\Sigma}^{-1/2} \hat{V} \hat{\Gamma}^{1/2} \hat{\Sigma}^{-1/2}$$

$$\hat{\nu} \leftarrow \max \left(1, \sqrt{\text{Tr}(\hat{\Sigma})} + 2\sqrt{2\|\hat{\Sigma}\|_{\text{op}} \log n} \right), \quad p \leftarrow \max(2, \log(1/\delta)), \quad \hat{\gamma} \leftarrow \frac{1}{2} \hat{\nu}^{-1} p^{-1/2} \|\hat{\Sigma}\|_{\text{op}}^{1/2}$$

$$\hat{L}(A, B) \leftarrow n^{-1} \sum_{i=1}^n \left\| \hat{\gamma}^{-1} \|\tilde{B}\|_F \tilde{A} \frac{\sum_{j=1}^n \exp(\hat{\gamma} x_i^\top \frac{\tilde{B}}{\|\tilde{B}\|_F} x_j) x_j}{\sum_{j=1}^n \exp(x_i^\top B x_j)} - y_i \right\|_2^2$$

$$\hat{R}(\tilde{A}, \tilde{B}) \leftarrow \frac{1}{8} \left\| \hat{\Sigma}^{1/2} (\tilde{A}^\top \tilde{A} - \tilde{B}^\top \tilde{B}) \hat{\Sigma}^{1/2} \right\|_F^2$$

$$\hat{Q}(\tilde{A}, \tilde{B}) \leftarrow \tilde{L}(\tilde{A}, \tilde{B}) + \hat{R}(\tilde{A}, \tilde{B})$$

for $t = 1, \dots, m$ **do**

$$\begin{cases} \tilde{A}_t \leftarrow A_{t-1} - \eta \nabla_A \hat{Q}(\tilde{A}_{t-1}, \tilde{B}_{t-1}) \\ \tilde{B}_t \leftarrow B_{t-1} - \eta \hat{\Sigma}^{-1} \nabla_B \hat{Q}(\tilde{A}_{t-1}, \tilde{B}_{t-1}) \end{cases}$$

end
return $(\hat{\gamma}^{-1} \|\tilde{B}_m\|_F \tilde{A}_m, \hat{\gamma} \|\tilde{B}_m\|_F^{-1} \tilde{B}_m)$

empirical estimate $\hat{\Sigma}$ in the definition of $R(\theta)$ and \hat{P} is defined analogously. Intuitively, when the number of samples n grows large, one should expect that $\nabla \hat{L}(\theta) \approx \nabla L(\theta)$ and $\nabla \hat{R}(\theta) \approx \nabla R(\theta)$. Gradient descent on $\hat{L}(\theta) + \hat{R}(\theta)$ should thus converge to a point near \mathcal{S} , provided that the gradient descent algorithm is initialized sufficiently close to \mathcal{S} .

However, this seemingly-natural idea has a key problem. The softmax function will not concentrate around its mean if its arguments are too large. It is thus necessary to rescale the parameter B to ensure that it does not grow too large during the training procedure. Our algorithm thus sets $(A, B) = (\gamma^{-1} \|\tilde{B}\|_F \tilde{A}, \gamma \|\tilde{B}\|_F^{-1} \tilde{B})$ and runs gradient descent on (\tilde{A}, \tilde{B}) . The regularized \tilde{R} is chosen so that the regularized loss \tilde{Q} is one-point strongly convex in (\tilde{A}, \tilde{B}) .

The following theorem shows that the population loss approaches its optimal value exponentially fast in the number of gradient descent iterations, up to an additional error term that decreases as $1/n$ in the size of the training set.

Theorem 1 *Let $\{\theta_k\}_{k=0}^m$ denote the sequence of iterates generated by Algorithm 1. With probability $1 - \delta$, the following inequality holds for all sufficiently large n :*

$$L(\theta_m) - L^* \leq K_1 \frac{\log^5(n) \log(mn/\delta)}{n} + K_2 \exp(-K_3 m),$$

where $K_1, K_2, K_3 \geq 0$ are constants depending on M and Σ and independent of m, n, δ .

References

- [1] Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Linear attention is (maybe) all you need (to understand transformer optimization). *arXiv preprint arXiv:2310.01082*, 2023.
- [2] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [3] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [4] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *Advances in neural information processing systems*, 36, 2024.
- [5] Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv preprint arXiv:2402.19442*, 2024.
- [6] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [7] Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- [8] Vladimir Koltchinskii and Karim Lounici. Concentration inequalities and moment bounds for sample covariance operators. *Bernoulli*, pages 110–133, 2017.
- [9] Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: Towards a mechanistic understanding. In *International Conference on Machine Learning*, pages 19689–19729. PMLR, 2023.
- [10] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [11] Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via procrustes flow. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 964–973, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/tu16.html>.
- [12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- [13] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*, volume 47. Cambridge University Press, 2018.
- [14] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.

- [15] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024.

Appendix A. Appendix

A.1. Optimization

A.2. Proof of Theorem 1

Proof We assume that the “good events” described in Lemma 26 occur. These events imply that Algorithm 1 selects an initial set of parameters which is close to the manifold of global minima, and that the empirical quantities $\hat{\Sigma}$, $\hat{\gamma}$, etc. are close to their population counterparts. Lemma 26 shows that these events simultaneously occur with probability $1 - \delta/2$, provided that n is sufficiently large. Let $Q(\theta)$ and \mathcal{S} be as described in Theorem ???. Define the rescaling function

$$g((A, B), \lambda) = (\lambda^{-1}A, \lambda B)$$

for all $\lambda > 0$. Notice that for any fixed $\lambda_0 > 0$, $g(\theta, \lambda_0)$ is a linear and invertible function of θ . In particular,

$$g(g(\theta, \lambda_0), \lambda_0^{-1}) = \theta.$$

Let

$$\tilde{\mathcal{S}} = \{g(\theta, \lambda) \mid \theta \in \mathcal{S}, \lambda > 0\}.$$

Notice that $L(\theta) = L^*$ for all $\theta \in \tilde{\mathcal{S}}$ because the function $L(\theta) = L(g(\theta, \lambda))$ for all $\lambda > 0$.

Define the scale-invariant regularized population loss

$$\tilde{Q}(\theta) = \inf_{\lambda > 0} Q(g(\theta, \lambda)).$$

It is easy to check that the optimal value of λ is

$$\lambda^*(A, B) = \left(\frac{\|\Sigma^{1/2} A^\top A \Sigma^{1/2}\|_F}{\|\Sigma^{1/2} B^\top \Sigma B \Sigma^{1/2}\|_F} \right)^{1/4}.$$

Define the empirical regularized loss

$$\hat{Q}(\theta) = \hat{L}(\theta) + \hat{R}(g(\theta, \hat{\lambda}^*)),$$

where we define

$$\begin{aligned} \hat{L}(A, B) &= \frac{1}{n} \sum_{i=1}^n \left\| A \frac{\sum_{j=1}^n \exp(x_i^\top B x_j) x_j}{\sum_{j=1}^n \exp(x_i^\top B x_j)} - y_i \right\|_2^2, \\ \hat{R}(A, B) &= \frac{1}{8} \left\| \hat{\Sigma}^{1/2} (A^\top A - B^\top \hat{\Sigma} B) \hat{\Sigma}^{1/2} \right\|_F^2, \\ \hat{\lambda}^*(A, B) &= \left(\frac{\|\hat{\Sigma}^{1/2} A^\top A \hat{\Sigma}^{1/2}\|_F}{\|\hat{\Sigma}^{1/2} B^\top \hat{\Sigma} B \hat{\Sigma}^{1/2}\|_F} \right)^{1/4}. \end{aligned}$$

Notice that $\lambda^*(g(\theta, \lambda)) = \lambda^*(\theta)$ for any $\lambda > 0$. Let us introduce the abbreviated notation $\hat{\lambda}_t^* = \hat{\lambda}^*(\theta_t)$ and $\bar{\theta}_t = g(\theta_t, \hat{\lambda}_t^*)$. Define

$$\hat{Z}_t = \left(\hat{\lambda}_t^* \nabla_A \hat{Q}(\theta_t), \hat{\lambda}_t^{-*} \hat{\Sigma}^{-1} \nabla_B \hat{Q}(\theta_t) \right).$$

Notice that the gradient update used by Algorithm 1 is exactly

$$\tilde{\theta}_{t+1} = \bar{\theta}_t - \eta \hat{Z}_t, \quad \theta_{t+1} = g(\tilde{\theta}_t, \hat{\gamma} \|\tilde{B}_t\|_F^{-1}).$$

In other words, in each iteration Algorithm 1 first preforms a gradient step and then a renormalization step. Let θ_t^* denote a projection in P -norm of the point $\bar{\theta}_t$ onto \mathcal{S} . Define the residual

$$\Delta_t = \bar{\theta}_t - \theta_t^*$$

and the potential

$$\phi_t = \|\Delta_t\|_P^2.$$

Notice that on the events described in Lemma 26, $\phi_0 \leq \varepsilon_0$. We will show that in each iteration, ϕ_t shrinks by a multiplicative factor, up to a small additive error. This will show that the sequence of iterates $\{\bar{\theta}_t\}_{t=1}^m$ converges (approximately) to $\tilde{\mathcal{S}}$. Since $L(\theta_t) = L(\bar{\theta}_t)$, this will show that $L(\theta_m)$ is near the optimal loss L^* . It is clear that

$$\phi_{t+1} \leq \|\bar{\theta}_{t+1} - \theta_t^*\|_P^2,$$

because the point θ_t^* can only be farther away from $\bar{\theta}_{t+1}$ than θ_{t+1}^* . Let \hat{Q} , Q^{emp} , and \tilde{Q}^{emp} be defined as in Section B. We see that

$$\begin{aligned} \phi_{t+1} &\leq \left\| \bar{\theta}_t - \eta \hat{Z}_t - \theta_t^* \right\|_P^2 \\ &= \|\Delta_t\|_P^2 + \eta^2 \|\hat{Z}_t\|_P^2 - 2\eta \langle \hat{Z}_t, \Delta_t \rangle_P \\ &= \mu(\eta) \|\Delta_t\|_P^2 + \nu_1 \|\zeta_t^{\text{emp}}\|_F^2 + \nu_2 \|\hat{\zeta}_t\|_F^2, \end{aligned}$$

where we defined the residuals

$$\hat{\zeta}_t = \nabla \hat{Q}(\theta_t) - \nabla \tilde{Q}^{\text{emp}}(\theta_t), \quad \zeta_t^{\text{emp}} = \nabla Q^{\text{emp}}(\bar{\theta}_t) - \nabla Q(\bar{\theta}_t),$$

and set

$$\mu(\nu) = (1 - 2\tilde{\alpha}\eta + \tilde{\beta}\eta^2)$$

and set

$$\nu_1 = 2a_1\eta + b_1\eta^2, \quad \nu_2 = 2a_2\eta + b_2\eta^2$$

where $\tilde{\alpha}, \tilde{\beta}, a_1, a_2, b_1, b_2$ are defined as in Lemma 4. Optimizing over η , we see that $\mu(\eta)$ is minimized when $\eta = \tilde{\alpha}\tilde{\beta}^{-1}$, in which case $\mu(\eta) = 1 - \tilde{\alpha}^2\tilde{\beta}^{-1}$. We emphasize that $0 < \mu < 1$. Theorem 9 shows that the variables $\{\zeta_t^{\text{emp}}, \hat{\zeta}_t\}_{t=0}^{m-1}$ are uniformly bounded with probability $1 - \delta/2$:

$$\sup_{t \in [m]} \max \left\{ \|\zeta_{t-1}^{\text{emp}}\|_F^2, \|\hat{\zeta}_{t-1}\|_F^2 \right\} \leq \frac{K \log^5(n) \log(mn/\delta)}{n}$$

for some constant $K \geq 0$. A standard inductive argument shows that

$$\phi_m \leq \frac{(\nu_1 + \nu_2)K \log^5(n) \log(mn/\delta)}{(1 - \mu)n} + \varepsilon_0 \mu^m.$$

Applying Lemma 8, we see that this implies that

$$Q(\bar{\theta}_m) - Q^* \leq \frac{K_1 \log^5(n) \log(mn/\delta)}{n} + K_2 \mu^m$$

for some $K_1, K_2 \geq 0$, where we set $Q^* = \inf_{\theta} Q(\theta)$. Observing that $Q(\bar{\theta}_m) \geq L(\bar{\theta}_m) = L(\theta_m)$ and $Q^* = L^*$ finishes the proof. ■

A.3. Proof of Theorem ??

For the readers convenience, we restate the theorem, this time showing the constants $\alpha, \beta, \varepsilon_0$ explicitly.

Theorem *The population loss $L(\theta)$ and the regularized population loss $Q(\theta)$ have the following properties:*

1. *The population loss can be written as*

$$L(\theta) = L^* + \frac{1}{2} \left\| A \Sigma B^\top \Sigma^{1/2} - M \Sigma^{1/2} \right\|_F^2,$$

where $L^* = \frac{1}{2} \text{Tr}(\Omega)$ is the irreducible loss.

2. *Define the regularized population loss*

$$Q(\theta) = L(\theta) + R(\theta),$$

where we set

$$R(\theta) = \frac{1}{8} \left\| \Sigma^{1/2} (A^\top A - B^\top \Sigma B) \Sigma^{1/2} \right\|_F^2.$$

Let $U \Gamma V^\top$ be a singular value decomposition of $M \Sigma^{1/2}$, where $U \in \mathbb{R}^{p \times d}$ and $V \in \mathbb{R}^{d \times d}$ satisfy $U^\top U = V^\top V = I_d$ and $\Gamma \in \mathbb{R}^{d \times d}$ is diagonal. Let \mathcal{S} be the smooth manifold consisting of points of the form

$$\begin{bmatrix} A \\ B \end{bmatrix} = \begin{bmatrix} U \Gamma^{1/2} J^\top \Sigma^{-1/2} \\ \Sigma^{-1/2} V \Gamma^{1/2} J^\top \Sigma^{-1/2} \end{bmatrix}$$

for some $J \in \mathbb{O}_d$. Each point $\theta \in \mathcal{S}$ is a global minimum of $Q(\cdot)$, and in particular satisfies $L(\theta) = L^*$ and $R(\theta) = 0$.

3. *Define the $(p+d) \times (p+d)$ block-diagonal matrix $P = \text{diag}(I_p, \Sigma)$. Define the P -weighted inner product*

$$\langle \theta_1, \theta_2 \rangle_P = \text{Tr}(\theta_1^\top P \theta_2)$$

and the associated P -norm

$$\|\theta\|_P = \sqrt{\langle \theta, \theta \rangle_P}.$$

The regularized population loss $Q(\theta)$ exhibits the following ‘‘one-point strong convexity’’ and ‘‘one-point smoothness’’ properties. Let θ^* denote a projection in the P -norm of θ onto \mathcal{S} . Let

$$\varepsilon_0 = \min \left\{ 1, \frac{K_0}{12K_1 \|\Sigma\|_{\text{op}}}, \sqrt{\frac{K_0}{8\|\Sigma\|_{\text{op}}^2}}, \frac{\sqrt{\|M \Sigma^{1/2}\|_*}}{\sqrt{\|\Sigma\|_{\text{op}}}}, \frac{\|M \Sigma^{1/2}\|_F}{6\sqrt{\|\Sigma\|_{\text{op}} \|M \Sigma^{1/2}\|}} \right\}.$$

For all θ which are ε_0 -close to \mathcal{S} in the P -norm, the following bounds hold:

$$\alpha \|\theta - \theta^*\|_P^2 \leq \langle P^{-1} \nabla Q(\theta), \theta - \theta^* \rangle_P, \quad (\text{one-point strong convexity}) \quad (3)$$

$$\|P^{-1} \nabla Q(\theta)\|_P^2 \leq \beta \|\theta - \theta^*\|_P^2, \quad (\text{one-point strong smoothness}) \quad (4)$$

where we define

$$\alpha = \frac{K_0}{4}, \quad \beta = \left(\frac{7}{2} + \frac{7}{4} \kappa^2(\Sigma) \right) K_1^2 + \frac{21}{4} \|\Sigma\|_{\text{op}}^2 K_1 + \frac{7}{4} \|\Sigma\|_{\text{op}}^4,$$

and K_0, K_1 are defined as in Lemma 2.

Proof To prove the first part of Theorem ??, we first recall that the population loss is

$$L(\theta) = \frac{1}{2} \mathbb{E}_{x_1, z_1} \left\| A \left(\frac{\mathbb{E}_{x_2}[\exp(x_1^\top B x_2) x_2]}{\mathbb{E}_{x_2}[\exp(x_1^\top B x_2)]} \right) - (M x_1 + z_1) \right\|_2^2.$$

Each of the expectations over x_2 can be computed using a standard completion-of-squares argument:

$$\mathbb{E}_{x_2}[\exp(x_1^\top B x_2) x_2] = \exp\left(\frac{1}{2} x_1^\top B \Sigma B^\top x_1\right) \Sigma B^\top x_1, \quad \mathbb{E}_{x_2}[\exp(x_1^\top B x_2)] = \exp\left(\frac{1}{2} x_1^\top B \Sigma B^\top x_1\right).$$

Canceling terms, we see that $L(\theta)$ can be written as

$$L(\theta) = \frac{1}{2} \mathbb{E}_{x_1, z_1} \left\| A \Sigma B^\top x_1 - (M x_1 + z_1) \right\|_2^2.$$

Integrating with respect to x_1 and z_1 and using the fact that $\mathbb{E}[x_1 z_1^\top] = 0$ because x_1 and z_1 are independent, we see that

$$L(\theta) = \frac{1}{2} \text{Tr}(\Omega) + \frac{1}{2} \left\| A \Sigma B^\top \Sigma^{1/2} - M \Sigma^{1/2} \right\|_F^2.$$

We now turn to the second part of Theorem ?. Recall that $Q(\theta) = L(\theta) + R(\theta)$, where $L(\theta) \geq L^*$ and $R(\theta) \geq 0$. It follows that any point $\theta = (A, B)$ such that $L(\theta) = L^*$ and $R(\theta) = 0$ is a global minimizer of $Q(\theta)$. The condition $L(\theta) = L^*$ implies that

$$A \Sigma B^\top \Sigma^{1/2} = M \Sigma^{1/2},$$

while the condition $R(\theta) = 0$ implies that

$$A^\top A = B^\top \Sigma B.$$

It is easy to check that all pairs (A, B) of the form

$$\begin{bmatrix} A \\ B \end{bmatrix} = \begin{bmatrix} U \Gamma^{1/2} J^\top \Sigma^{-1/2} \\ \Sigma^{-1/2} V \Gamma^{1/2} J^\top \Sigma^{-1/2} \end{bmatrix}$$

satisfy both equations. We now prove the third part of Theorem ?. The proof depends crucially on the following lemma. The significance of this lemma is that it shows that the first-order condition satisfied by θ^* implies a symmetry condition which allows us to establish strong convexity of $Q(\theta)$. This symmetry condition is the key reason why we choose the P -norm to measure the distance between θ and \mathcal{S} .

Lemma 2 Fix any $\theta = (A, B)$ where $A \in \mathbb{R}^{p \times d}$ and $B \in \mathbb{R}^{d \times d}$ and let θ^* be the projection of θ onto \mathcal{S} in the P -norm. The matrix $\Delta^\top P \theta^* \Sigma$ is symmetric.

We defer the proof of Lemma 2 to the subsequent section.

We first establish one-point strong convexity before proving one-point smoothness. Our proof is inspired by the proof of Theorem 3.3 in [11], but is modified to account for the fact that $L(\theta)$ is not symmetric in A and B due to the extra factor of $\Sigma^{1/2}$ attached to B . An algebraic calculation shows that $Q(\theta)$ can be rewritten in the form

$$Q(\theta) = L^* + \frac{1}{8} \left\| P^{1/2} (\theta \Sigma \theta^\top - 2 \text{Sym}(M)) P^{1/2} \right\|_F^2 - \frac{1}{2} \|M \Sigma^{1/2}\|_F^2,$$

where we define

$$\text{Sym}(M) = \begin{bmatrix} 0 & M \\ M^\top & 0 \end{bmatrix}.$$

Define

$$\tilde{\theta}^* = \begin{bmatrix} I_p & 0 \\ 0 & -I_d \end{bmatrix} \theta^*.$$

Notice that

$$2\text{Sym}(M) = \theta^* \Sigma \theta^{*\top} - \tilde{\theta}^* \Sigma \tilde{\theta}^{*\top} \quad (5)$$

and

$$\theta^{*\top} P \tilde{\theta}^* = 0. \quad (6)$$

Set $\Delta = \theta - \theta^*$. Applying (5), we see that

$$\begin{aligned} \langle P^{-1} \nabla Q(\theta), \Delta \rangle_P &= \langle \nabla Q(\theta), \Delta \rangle \\ &= \frac{1}{2} \langle P(\theta \Sigma \theta^\top - 2\text{Sym}(M)) P \theta \Sigma, \Delta \rangle \\ &= \frac{1}{2} \langle P(\theta \Sigma \theta^\top - \theta^* \Sigma \theta^{*\top}) P \theta \Sigma, \Delta \rangle + \frac{1}{2} \langle P \tilde{\theta}^* \Sigma \tilde{\theta}^{*\top} P \theta \Sigma, \Delta \rangle. \end{aligned} \quad (7)$$

We lower-bound each term of (7) separately. It is convenient to lower-bound the second term of (7) first. Notice that

$$\frac{1}{2} \langle P \tilde{\theta}^* \Sigma \tilde{\theta}^{*\top} P \theta \Sigma, \Delta \rangle = \frac{1}{2} \text{Tr}(\theta^\top P \tilde{\theta}^* \Sigma \tilde{\theta}^{*\top} P \theta \Sigma) - \frac{1}{2} \text{Tr}(\theta^{*\top} P \tilde{\theta}^* \Sigma \tilde{\theta}^{*\top} P \theta \Sigma).$$

The first term is non-negative because it is the trace of the product of two psd matrices. We see that the second term is equal to zero in light of (6). This proves that the second term of (7) is non-negative.

We now show that the first term of (7) is bounded below by a constant multiple of $\|\Delta\|_F^2$, provided that θ is sufficiently close to \mathcal{S} . We observe that

$$\begin{aligned} \frac{1}{2} \langle P(\theta \Sigma \theta^\top - \theta^* \Sigma \theta^{*\top}) P \theta \Sigma, \Delta \rangle &= \frac{1}{2} \text{Tr} \left(\Delta^\top P(\theta \Sigma \theta^\top - \theta^* \Sigma \theta^{*\top}) P \theta \Sigma \right) \\ &= \frac{1}{2} \text{Tr} \left(\Delta^\top P(\theta^* \Sigma \Delta^\top + \Delta \Sigma \theta^{*\top} + \Delta \Sigma \Delta^\top) P(\theta^* + \Delta) \Sigma \right) \\ &= S + T, \end{aligned}$$

where we set S be the sum of terms which are quadratic in Δ and set T be the sum of all remaining terms:

$$\begin{aligned} S &= \frac{1}{2} \text{Tr} \left(\Delta^\top P \theta^* \Sigma \Delta^\top P \theta^* \Sigma + \Delta^\top P \Delta \Sigma \theta^{*\top} P \theta^* \Sigma \right), \\ T &= \frac{1}{2} \text{Tr} \left(\Delta^\top P \Delta \Sigma \Delta^\top P \theta^* \Sigma + \Delta^\top P \Delta \Sigma \Delta^\top P \Delta \Sigma + \Delta^\top P \theta^* \Sigma \Delta^\top P \Delta \Sigma + \Delta^\top P \Delta \Sigma \theta^{*\top} P \Delta \Sigma \right). \end{aligned}$$

We lower-bound S and T individually. We see that the first term of S is equal to $\frac{1}{2} \|\Delta^\top P \theta^* \Sigma\|_F^2$ using the symmetry condition described in Lemma 2, while the second term is equal to $\frac{1}{2} \|P^{1/2} \theta^* \Sigma \Delta^\top P^{1/2}\|_F^2$. Applying elementary properties of the Frobenius norm, the fact that $\|X\|_P = \|P^{1/2} X\|_F$ for all matrices X , and Lemma 3, we see that

$$\begin{aligned} S &\geq \frac{1}{2} \sigma_{\min}^2(P^{1/2} \theta^* \Sigma) \|\Delta\|_P^2 \\ &\geq \frac{1}{2} K_0 \|\Delta\|_P^2. \end{aligned}$$

Applying the Cauchy-Schwarz inequality, Lemma 2, and elementary properties of the Frobenius norm, we see that

$$T \geq -\frac{3}{2} K_1 \|\Sigma\|_{\text{op}} \|\Delta\|_P^3 - \|\Sigma\|_{\text{op}}^2 \|\Delta\|_P^4.$$

Putting the pieces together, we see that

$$\langle P^{-1}\nabla Q(\theta), \Delta \rangle_P \geq \frac{1}{4}K_0\|\Delta\|_P^2$$

provided that

$$\|\Delta\|_P \leq \min \left\{ \frac{K_0}{12K_1\|\Sigma\|_{\text{op}}}, \sqrt{\frac{K_0}{8\|\Sigma\|_{\text{op}}^2}} \right\}.$$

We now establish strong smoothness of $Q(\theta)$ near \mathcal{S} . We see that

$$\begin{aligned} \|P^{-1}\nabla Q(\theta)\|_P^2 &= \frac{1}{4}\|(\theta\Sigma\theta^\top - 2\text{Sym}(M))P\theta\Sigma\|_P^2 \\ &= \frac{1}{4}\|((\theta^* + \Delta)\Sigma(\theta^* + \Delta)^\top - 2\text{Sym}(M))P(\theta^* + \Delta)\Sigma\|_P^2 \\ &\leq \sum_{i=1}^7 T_i, \end{aligned}$$

where we define

$$\begin{aligned} T_1 &= \frac{7}{4}\|(\theta^*\Sigma\theta^{*\top} - 2\text{Sym}(M))P\Delta\Sigma\|_P^2, \\ T_2 &= \frac{7}{4}\|\theta^*\Sigma\Delta^\top P\theta^*\Sigma\|_P^2, \\ T_3 &= \frac{7}{4}\|\Delta\Sigma\theta^{*\top} P\theta^*\Sigma\|_P^2, \\ T_4 &= \frac{7}{4}\|\theta^*\Sigma\Delta^\top P\Delta\Sigma\|_P^2, \\ T_5 &= \frac{7}{4}\|\Delta\Sigma\theta^{*\top} P\Delta\Sigma\|_P^2, \\ T_6 &= \frac{7}{4}\|\Delta\Sigma\Delta^\top P\theta^*\Sigma\|_P^2, \\ T_7 &= \frac{7}{4}\|\Delta\Sigma\Delta^\top P\Delta\Sigma\|_P^2, \end{aligned}$$

and use the easily-verified algebraic fact that

$$(\theta^*\Sigma\theta^{*\top} - 2\text{Sym}(M))P\theta^*\Sigma = 0.$$

Recall that for every matrix X , $\|X\|_P = \|P^{1/2}X\|_F$. We recall that

$$\theta^*\Sigma\theta^{*\top} - 2\text{Sym}(M) = \tilde{\theta}^*\Sigma\tilde{\theta}^{*\top}$$

and observe that $\|P^{1/2}\tilde{\theta}^*\Sigma\|_{\text{op}} = \|P^{1/2}\theta^*\Sigma\|_{\text{op}}$. It is clear that $\|\Delta\|_P^4, \|\Delta\|_P^6 \leq \|\Delta\|_P^2$ for all Δ such that $\|\Delta\|_P^2 \leq 1$. Applying Lemma 3 and elementary properties of the Frobenius norm, we see that

$$\|P^{-1}\nabla Q(\theta)\|_P^2 \leq \left(\left(\frac{7}{2} + \frac{7}{4}\kappa^2(\Sigma) \right) K_1^2 + \frac{21}{4}\|\Sigma\|_{\text{op}}^2 K_1 + \frac{7}{4}\|\Sigma\|_{\text{op}}^4 \right) \|\Delta\|_P^2$$

for all Δ such that $\|\Delta\|_P \leq 1$. ■

A.4. Proofs of Lemma 2 and Lemma 3

We first prove Lemma 2:

Proof Recall from the second part of Theorem ?? that each point on \mathcal{S} has the form

$$\begin{bmatrix} U\Gamma^{1/2}J^\top\Sigma^{-1/2} \\ \Sigma^{-1/2}V\Gamma^{1/2}J^\top\Sigma^{-1/2} \end{bmatrix}$$

for some orthogonal $J \in \mathbb{R}^{d \times d}$, where UV^\top is the singular value decomposition of $M\Sigma^{1/2}$. It follows that

$$\theta^* = S\Gamma^{1/2}J^{*\top}\Sigma^{-1/2},$$

where we set

$$S = \begin{bmatrix} U \\ \Sigma^{-1/2}V \end{bmatrix}, \quad J^* = \operatorname{argmin}_{J \in \mathbb{O}_d} F(J)$$

and

$$F(J) = \frac{1}{2} \left\| \theta - S\Gamma^{1/2}J^\top\Sigma^{-1/2} \right\|_P^2.$$

We relax the constraint that J is orthogonal and define the Lagrangian

$$G(J, \Lambda) = F(J) + \frac{1}{2} \operatorname{Tr}(\Lambda(J^\top J - I_d)).$$

We may assume that $\Lambda \in \mathbb{R}^{d \times d}$ is symmetric without loss of generality, since the matrix $J^\top J - I_d$ is clearly symmetric and the trace of the product of a symmetric matrix and a skew-symmetric matrix is zero. Any orthogonal minimizer J^* of $F(J)$ satisfies the first-order condition $\nabla_J G(J, \Lambda) = 0$ for some symmetric Λ . We see that this condition is precisely

$$-\Sigma^{-1/2}\theta^\top P S \Gamma^{1/2} + \Sigma^{-1} J^* \Gamma^{1/2} (S^\top P S) \Gamma^{1/2} + J^* \Lambda = 0.$$

Observing that $S^\top P S = 2I_d$ and rearranging, we obtain

$$\Sigma^{-1/2}\theta^\top P S \Gamma^{1/2} = 2\Sigma^{-1} J^* \Gamma + J^* \Lambda. \quad (8)$$

Recall that $\Delta = \theta - \theta^*$ and $\theta^* = S\Gamma^{1/2}J^{*\top}\Sigma^{-1/2}$. It follows that

$$\begin{aligned} \Sigma^{-1/2}\theta^\top P S \Gamma^{1/2} &= \Sigma^{-1/2}\Delta^\top P S \Gamma^{1/2} + \Sigma^{-1/2}\theta^{*\top} P S \Gamma^{1/2} \\ &= \Sigma^{-1/2}\Delta^\top P S \Gamma^{1/2} + 2\Sigma^{-1} J^* \Gamma, \end{aligned}$$

where we once again used the identity $S^\top P S = 2I_d$. In light of (8), this equation implies that

$$\Sigma^{-1/2}\Delta^\top P S \Gamma^{1/2} = J^* \Lambda. \quad (9)$$

We use this identity to prove that $\Delta^\top P \theta^* \Sigma$ is symmetric. Plugging in our expression for θ^* once more, we see that

$$\begin{aligned} \Delta^\top P \theta^* \Sigma &= \Delta^\top P S \Gamma^{1/2} J^{*\top} \Sigma^{1/2} \\ &= \Sigma^{1/2} J^* \Lambda J^{*\top} \Sigma^{1/2}, \end{aligned}$$

where we applied (9) in the second step. This matrix is symmetric because Λ is symmetric. ■

Lemma 3 *Let $\theta^* = (A^*, B^*)$ be any point in \mathcal{S} . The following inequalities hold:*

$$K_0 \leq \sigma_{\min}^2(P^{1/2}\theta^*\Sigma), \quad \sigma_1^2(P^{1/2}\theta^*\Sigma) \leq K_1,$$

where we set

$$K_0 = 2\sigma_{\min}(M\Sigma^{1/2})\sigma_{\min}(\Sigma) \quad K_1 = 2\sigma_1(M\Sigma^{1/2})\sigma_1(\Sigma).$$

In addition, the following identity holds:

$$\Sigma^{1/2}A^{*\top}A^*\Sigma^{1/2} = \Sigma^{1/2}B^{*\top}\Sigma B^*\Sigma^{1/2} = \Gamma.$$

Proof Recall from the second part of Theorem ?? that each point $\theta^* = (A^*, B^*)$ on \mathcal{S} has the form

$$\begin{bmatrix} A^* \\ B^* \end{bmatrix} = \begin{bmatrix} U\Gamma^{1/2}J^\top\Sigma^{-1/2} \\ \Sigma^{-1/2}V\Gamma^{1/2}J^\top\Sigma^{-1/2} \end{bmatrix}$$

for some orthogonal $J \in \mathbb{R}^{d \times d}$, where UTV^\top is the singular value decomposition of $M\Sigma^{1/2}$. We see that

$$P^{1/2}\theta^*\Sigma = \begin{bmatrix} U \\ V \end{bmatrix} \Gamma^{1/2}J^\top\Sigma^{1/2}.$$

Using the fact that

$$U^\top U + V^\top V = 2I_d$$

and the fact that J^* is orthogonal, and applying elementary properties of singular values, we obtain the bound

$$2\sigma_d(M\Sigma^{1/2})\sigma_d(\Sigma) \leq \sigma_d^2(P^{1/2}\theta^*\Sigma).$$

The upper bound

$$\sigma_1^2(P^{1/2}\theta^*\Sigma) \leq 2\sigma_1(M\Sigma^{1/2})\sigma_1(\Sigma)$$

follows from a similar calculation. We observe via direct calculation that

$$\Sigma^{1/2}A^{*\top}A^*\Sigma^{1/2} = \Sigma^{1/2}B^{*\top}\Sigma B^*\Sigma^{1/2} = \Gamma.$$

■

A.5. Approximate strong convexity and strong smoothness of the empirical loss

The following lemma shows that the strong convexity and smoothness established in Theorem ?? for the population loss holds approximately on the empirical loss:

Lemma 4 *Let α, β be as described in Theorem ?? and let λ_{\max} be defined as in Lemma 7. Using the notation defined in the proof of Theorem 1, the following inequalities hold for each $t \in [m]$:*

$$\langle \hat{Z}_t, \Delta_t \rangle_P \geq \tilde{\alpha} \|\Delta_t\|_P^2 - a_1 \|\zeta_t^{\text{emp}}\|_F^2 - a_2 \|\hat{\zeta}_t\|_F^2 \tag{10}$$

$$\|\hat{Z}_t\|_P^2 \leq \tilde{\beta} \|\Delta\|_P^2 + b_1 \|\zeta_t^{\text{emp}}\|_F^2 + b_2 \|\hat{\zeta}_t\|_F^2, \tag{11}$$

where we define the constants

$$\begin{aligned}\tilde{\alpha} &= \frac{\alpha}{2}, \\ a_1 &= 9\alpha^{-1}, \\ a_2 &= 64\alpha^{-1}\lambda_{\max}^2, \\ \tilde{\beta} &= 3\beta\|P\|_{\text{op}}^2 + \frac{3\alpha^2}{32}, \\ b_1 &= 3\|P\|_{\text{op}} + 12, \\ b_2 &= 96\lambda_{\max}^2\end{aligned}$$

and the random variables

$$\begin{aligned}\hat{\zeta}_t &= \nabla\hat{Q}(\theta_t) - \nabla\tilde{Q}^{\text{emp}}(\theta_t), \\ \zeta_t^{\text{emp}} &= \nabla Q^{\text{emp}}(\bar{\theta}_t) - \nabla Q(\bar{\theta}_t),\end{aligned}$$

where \hat{Q} , Q^{emp} , and \tilde{Q}^{emp} are formally defined in Section B.

Proof Let Q^{emp} be the function obtained by replacing the true covariance Σ by the empirical covariance $\hat{\Sigma}$ in the definition of Q :

$$Q^{\text{emp}}(A, B) = \|A\hat{\Sigma}B\|_F^2 + \|\hat{\Sigma}^{1/2}(A^\top A - B^\top \hat{\Sigma}B)\hat{\Sigma}^{1/2}\|_F^2.$$

Notice that the multiplier $\hat{\lambda}_t^*$ is precisely the choice of λ which minimizes $Q^{\text{emp}}(g(\theta_t), \lambda)$. Lastly, define $\tilde{Q}^{\text{emp}}(\theta) = Q^{\text{emp}}(g(\theta), \hat{\lambda}^*(\theta))$. Define the error residual

$$\xi_t = \hat{Z}_t - \left(\hat{\lambda}_t^* \nabla_A \tilde{Q}^{\text{emp}}(\theta_t), \hat{\lambda}_t^{-*} \nabla_B \Sigma^{-1} \tilde{Q}^{\text{emp}}(\theta_t) \right).$$

We first establish the approximate one-point strong convexity property (10). We observe that

$$\begin{aligned}\langle \hat{Z}_t, \Delta_t \rangle_P &= \left\langle \left(\hat{\lambda}_t^* \nabla_A \tilde{Q}^{\text{emp}}(\theta_t), \hat{\lambda}_t^{-*} \nabla_B \tilde{Q}^{\text{emp}}(\theta_t) \right) + P \xi_t, \Delta_t \right\rangle \\ &= \langle \nabla Q^{\text{emp}}(\bar{\theta}_t) + P \xi_t, \Delta_t \rangle \\ &= \langle \nabla Q(\bar{\theta}_t), \Delta_t \rangle + \langle \zeta_t^{\text{emp}} + P \xi_t, \Delta_t \rangle \\ &= \langle P^{-1} \nabla Q(\bar{\theta}_t), \Delta_t \rangle_P + \langle \zeta_t^{\text{emp}} + P \xi_t, \Delta_t \rangle \\ &\geq \alpha \|\Delta_t\|_P^2 + \langle \zeta_t^{\text{emp}} + P \xi_t, \Delta_t \rangle,\end{aligned}$$

where we applied Lemma 6 and the one-point strong convexity property established in Theorem ???. We now lower-bound

$$\langle \zeta_t^{\text{emp}} + P \xi_t, \Delta_t \rangle.$$

Applying the Cauchy-Schwarz inequality, we see that

$$\langle \zeta_t^{\text{emp}}, \Delta_t \rangle \geq -\|\zeta_t^{\text{emp}}\|_F \|\Delta_t\|_F.$$

We recall Young's inequality: when $u, v \geq 0$ then $uv \leq \frac{1}{2\nu}u^2 + \frac{\nu}{2}v^2$ for any $\nu > 0$. Applying this inequality, we see that

$$\langle \zeta_t^{\text{emp}}, \Delta_t \rangle \geq -\alpha^{-1} \|\zeta_t^{\text{emp}}\|_F^2 - \frac{\alpha}{4} \|\Delta_t\|_F^2.$$

Similarly, we see that

$$\langle P\xi_t, \Delta_t \rangle \geq -\|P\xi_t\|_F \|\Delta_t\|_F,$$

which is lower bounded by

$$-2\alpha^{-1}\|P\xi_t\|_F^2 - \frac{\alpha}{8}\|\Delta_t\|_F^2.$$

Applying Lemma 5 to bound the first term yields the claim.

We now establish the approximate one-point smoothness property (11). We observe that

$$\begin{aligned} \|\hat{Z}_t\|_P^2 &= \left\| \left(\hat{\lambda}_t^* \nabla_A \tilde{Q}^{\text{emp}}(\theta_t), \hat{\lambda}_t^{-*} \nabla_B \tilde{Q}^{\text{emp}}(\theta_t) \right) + P\xi_t \right\|_P^2 \\ &= \|\nabla Q^{\text{emp}}(\bar{\theta}_t) + P\xi_t\|_P^2 \\ &\leq 3\|\nabla Q(\bar{\theta}_t)\|_P^2 + 3\|\zeta_t^{\text{emp}}\|_P^2 + 3\|P\xi_t\|_P^2 \\ &\leq 3\|P\|_{\text{op}}^2\|P^{-1}\nabla Q(\bar{\theta}_t)\|_P^2 + 3\|P\|_{\text{op}}\|\zeta_t^{\text{emp}}\|_F^2 + 3\|P\xi_t\|_P^2 \\ &\leq 3\beta\|P\|_{\text{op}}^2\|\Delta_t\|_P^2 + 3\|P\|_{\text{op}}\|\zeta_t^{\text{emp}}\|_F^2 + 3\|P\xi_t\|_P^2, \end{aligned}$$

where applied Lemma 6 and the one-point strong smoothness property established in Theorem ???. Applying Lemma 5 to bound the last term yields the claim. \blacksquare

Lemma 5 *Assume that the events described in Lemma 26 occur. Let λ_{\max} be defined as in Lemma 7. Using the notation defined in the proofs of Theorem 1 and Lemma 4, the following inequality holds:*

$$\|P\xi_t\|_F^2 \leq \frac{\alpha^2}{16}\|\Delta_t\|_P^2 + 4\|\zeta_t^{\text{emp}}\|_F^2 + 32\lambda_{\max}^2\|\hat{\zeta}_t\|_F^2.$$

Proof We see that

$$P\xi_t = \left(\hat{\lambda}_t \hat{\zeta}_t^A, \hat{\lambda}_t^{-*} \Sigma \hat{\Sigma}^{-1} \hat{\zeta}_t^B \right) + \left(0, \hat{\lambda}_t^{-*} (\Sigma \hat{\Sigma}^{-1} - I_d) \nabla_B \tilde{Q}^{\text{emp}}(\theta_t) \right),$$

where $\hat{\zeta}_t^A$ and $\hat{\zeta}_t^B$ denote the A and B components of $\hat{\zeta}_t$, respectively. It follows that

$$\|P\xi_t\|_F^2 \leq 2 \left\| \left(\hat{\lambda}_t \hat{\zeta}_t^A, \hat{\lambda}_t^{-*} \Sigma \hat{\Sigma}^{-1} \hat{\zeta}_t^B \right) \right\|_F^2 + 2 \left\| \hat{\lambda}_t^{-*} (\Sigma \hat{\Sigma}^{-1} - I_d) \nabla_B \tilde{Q}^{\text{emp}}(\theta_t) \right\|_F^2. \quad (12)$$

We bound each of these terms in turn. The first term of (12) is at most

$$2 \max\{\hat{\lambda}_t^2, \hat{\lambda}_t^{-2}\} \max\{1, \|\Sigma \hat{\Sigma}^{-1}\|_{\text{op}}^2\} \|\hat{\zeta}_t\|_F^2.$$

We observe that

$$\begin{aligned} \max\{\hat{\lambda}_t, \hat{\lambda}_t^{-*}\} &\leq \max\{\lambda_t^*, \lambda_t^{-*}\} + \max\{|\hat{\lambda}_t^* - \lambda_t^*|, |\hat{\lambda}_t^{-*} - \lambda_t^{-*}|\} \\ &\leq \max\{\lambda_t^*, \lambda_t^{-*}\} + \varepsilon_\lambda \\ &\leq 2\lambda_{\max}, \end{aligned}$$

where we applied Lemma 7 and the fact that $\varepsilon_\lambda \leq \lambda_{\max}$. We observe that our assumption that $\varepsilon_\Sigma \leq \|\Sigma\|_{\text{op}}^{-1}$ implies that

$$\|\Sigma \hat{\Sigma}^{-1}\|_{\text{op}} \leq 2.$$

Putting the pieces together, we see that the first term of (12) is at most

$$32\lambda_{\max}^2\|\hat{\zeta}_t\|_F^2.$$

Applying Lemma 6, we see that the second term of (12) is equal to

$$2 \left\| (\Sigma \hat{\Sigma}^{-1} - I_d) \nabla_B Q^{\text{emp}}(\bar{\theta}_t) \right\|_F^2,$$

which is at most

$$4 \left\| (\Sigma \hat{\Sigma}^{-1} - I_d) \nabla_B Q(\bar{\theta}_t) \right\|_F^2 + 4 \left\| (\Sigma \hat{\Sigma}^{-1} - I_d) (\nabla_B Q^{\text{emp}}(\bar{\theta}_t) - \nabla_B Q(\bar{\theta}_t)) \right\|_F^2. \quad (13)$$

We bound each of these two terms separately. The first term of (13) is at most

$$4 \left\| (\Sigma \hat{\Sigma}^{-1} - I_d) \Sigma^{-1/2} \right\|_{\text{op}}^2 \left\| \Sigma^{1/2} \nabla_B Q(\bar{\theta}_t) \right\|_F^2.$$

Using the definition of the P -norm, we see that this is at most

$$4 \left\| (\Sigma \hat{\Sigma}^{-1} - I_d) \Sigma^{-1/2} \right\|_{\text{op}}^2 \left\| \nabla Q(\bar{\theta}_t) \right\|_P^2,$$

which in turn is at most

$$\frac{\alpha^2}{16} \left\| \Delta_t \right\|_P^2,$$

where applied Theorem ?? and used the assumption that

$$\varepsilon_\Sigma^2 \leq \frac{\alpha^2}{64\beta} \left\| \Sigma \right\|_{\text{op}}^{-1}.$$

We now bound the second term of (13). We see that this term is at most

$$4 \left\| \zeta_t^{\text{emp}} \right\|_F^2,$$

where we used the assumption that $\varepsilon_\Sigma \leq \left\| \Sigma \right\|_{\text{op}}^{-1}$. ■

Lemma 6 *Using the notation defined in the proofs of Theorem 1 and Lemma 4, the following identity holds:*

$$\nabla \tilde{Q}^{\text{emp}}(\theta_t) = \left(\hat{\lambda}_t^{-*} \nabla_A Q^{\text{emp}}(\bar{\theta}_t), \hat{\lambda}_t^* \nabla_B Q^{\text{emp}}(\bar{\theta}_t) \right).$$

Proof Let (δ_A, δ_B) be arbitrary direction in parameter space. We see that

$$\left\langle \nabla \tilde{Q}^{\text{emp}}(\theta_t), (\delta_A, \delta_B) \right\rangle$$

is equal to the sum of two terms, namely

$$\left\langle \nabla_A Q^{\text{emp}}(\bar{\theta}_t), \hat{\lambda}_t^{-*} \delta_A - \hat{\lambda}_t^{-2*} \langle \nabla \hat{\lambda}_t^*, (\delta_A, \delta_B) \rangle A \right\rangle$$

and

$$\left\langle \nabla_B Q^{\text{emp}}(\bar{\theta}_t), \hat{\lambda}_t^* \delta_B + \langle \nabla \hat{\lambda}_t^*, (\delta_A, \delta_B) \rangle B \right\rangle,$$

where we used the chain rule to account for the fact that $\hat{\lambda}^*(\theta)$ varies as θ varies. We see that this is equal to

$$\left\langle \nabla_A Q^{\text{emp}}(\bar{\theta}_t), \hat{\lambda}_t^{-*} \delta_A \right\rangle + \left\langle \nabla_B Q^{\text{emp}}(\bar{\theta}_t), \hat{\lambda}_t^* \delta_B \right\rangle$$

because $\hat{\lambda}_t^*$ satisfies the first-order condition

$$\left\langle \nabla_A Q^{\text{emp}}(\bar{\theta}_t), \hat{\lambda}_t^{-2*} A \right\rangle - \left\langle \nabla_B Q^{\text{emp}}(\bar{\theta}_t), B \right\rangle = 0$$

because $\hat{\lambda}_t^*$ minimizes the function $Q^{\text{emp}}(\lambda^{-1} A_t, \lambda B_t)$. The fact that this identity holds for arbitrary (δ_A, δ_B) implies the claim. ■

Lemma 7 *Let \mathcal{S} and ε_0 be defined as in Theorem ???. There exists a constant λ_{\max} such that for all θ which are ε_0 -close to \mathcal{S} , the following inequality holds:*

$$\max\{\lambda^*(\theta), \lambda^{-*}(\theta)\} \leq \lambda_{\max}.$$

Proof Let $\bar{\mathcal{S}}$ denote the ε_0 -neighborhood of \mathcal{S} . Recall that

$$\varepsilon_0 \leq \min \left\{ \frac{\sqrt{\|M\Sigma^{1/2}\|_*}}{\sqrt{\|\Sigma\|_{\text{op}}}}, \frac{\|M\Sigma^{1/2}\|_F}{6\sqrt{\|\Sigma\|_{\text{op}}\|M\Sigma^{1/2}\|_*}} \right\}.$$

We claim that

$$\sup_{(A,B) \in \bar{\mathcal{S}}} \max\{\lambda^*(A, B), \lambda^{-*}(A, B)\} \leq 3^{1/4}.$$

Indeed, fix any $(A, B) \in \bar{\mathcal{S}}$, and choose $(A^*, B^*) \in \mathcal{S}$ such that

$$\|A - A^*\|_F \leq \varepsilon_0, \quad \|\Sigma^{1/2}(B - B^*)\|_F \leq \varepsilon_0.$$

Using the identities described in Lemma 3, together with the triangle inequality, we see that

$$\left| \|\Sigma^{1/2}A^\top A\Sigma^{1/2}\|_F - \|M\Sigma^{1/2}\|_F \right| \leq \left(2\sqrt{\|M\Sigma^{1/2}\|_*} + \sqrt{\|\Sigma\|_{\text{op}}\varepsilon_0} \right) \sqrt{\|\Sigma\|_{\text{op}}\varepsilon_0},$$

which is at most

$$3\sqrt{\|\Sigma\|_{\text{op}}\|M\Sigma^{1/2}\|_*\varepsilon_0} \leq \frac{1}{2}\|M\Sigma^{1/2}\|_F.$$

Therefore

$$\frac{1}{2}\|M\Sigma^{1/2}\|_F \leq \|\Sigma^{1/2}A^\top A\Sigma^{1/2}\|_F \leq \frac{3}{2}\|M\Sigma^{1/2}\|_F.$$

The same argument gives

$$\frac{1}{2}\|M\Sigma^{1/2}\|_F \leq \|\Sigma^{1/2}B^\top \Sigma B\Sigma^{1/2}\|_F \leq \frac{3}{2}\|M\Sigma^{1/2}\|_F.$$

Hence

$$\lambda^*(A, B) = \left(\frac{\|\Sigma^{1/2}A^\top A\Sigma^{1/2}\|_F}{\|\Sigma^{1/2}B^\top \Sigma B\Sigma^{1/2}\|_F} \right)^{1/4} \leq 3^{1/4},$$

and similarly $\lambda^{-*}(A, B) \leq 3^{1/4}$. ■

Lemma 8 (Descent Lemma) *Let β , ε_0 and \mathcal{S} be defined as in Theorem ???. Suppose that θ is ε_0 -close to \mathcal{S} in the P -norm. Let θ^* denote the projection of θ onto \mathcal{S} in the P -norm and let $\Delta = \theta - \theta^*$. The following inequality holds:*

$$Q(\theta) - Q^* \leq \frac{\sqrt{\beta}}{2} \|\Delta\|_P^2.$$

Proof For all $t \in [0, 1]$, let $s(t) = \theta^* + t(\theta - \theta^*)$. Notice that

$$\begin{aligned}
 Q(\theta) - Q^* &= Q(\theta) - Q(\theta^*) \\
 &= Q(s(1)) - Q(s(0)) \\
 &= \int_0^1 \frac{d}{dt} Q(s(t)) dt \\
 &= \int_0^1 \langle \nabla Q(s(t)), \theta - \theta^* \rangle dt \\
 &= \int_0^1 \langle P^{-1} \nabla Q(s(t)), \theta - \theta^* \rangle_P dt \\
 &\leq \int_0^1 \|P^{-1} \nabla Q(s(t))\|_P \|\theta - \theta^*\|_P dt \\
 &\leq \int_0^1 \sqrt{\beta} t \|\theta - \theta^*\|_P^2 dt \\
 &= \frac{\sqrt{\beta}}{2} \|\Delta\|_P^2,
 \end{aligned}$$

where we applied the Cauchy-Schwartz inequality and the one-point smoothness property described in Theorem ??.

Appendix B. Concentration

Recall that our analysis is based on controlling the pairwise distances between the gradients of the five regularized loss functions $\hat{Q}(\theta)$, $Q(\theta)$, $Q^{\text{emp}}(\theta)$, $\tilde{Q}(\theta)$, and $\tilde{Q}^{\text{emp}}(\theta)$. In particular, we wish to argue that gradient descent on \hat{Q} behaves almost identically to gradient descent on $\tilde{Q}(\theta)$ when n is large. These functions are defined as follows:

$$\begin{aligned}
 \hat{Q}(\theta) &= \hat{L}(\theta) + \inf_{\lambda > 0} R^{\text{emp}}(g(\theta, \lambda)), \\
 Q(\theta) &= L(\theta) + R(\theta), \\
 \tilde{Q}(\theta) &= \inf_{\lambda > 0} Q(g(\theta, \lambda)), \\
 Q^{\text{emp}} &= L^{\text{emp}}(\theta) + R^{\text{emp}}(\theta), \\
 \tilde{Q}^{\text{emp}}(\theta) &= \inf_{\lambda > 0} Q^{\text{emp}}(g(\theta, \lambda)),
 \end{aligned}$$

where we define the rescaling function

$$g((A, B), \lambda) = (\lambda^{-1}A, \lambda B)$$

and set

$$\begin{aligned}
 L(\theta) &= \frac{1}{2} \|A\Sigma B^\top - M\Sigma^{1/2}\|_F^2, \\
 L^{\text{emp}}(\theta) &= \frac{1}{2} \|A\hat{\Sigma} B^\top - M\hat{\Sigma}^{1/2}\|_F^2, \\
 \hat{L} &= \frac{1}{2n} \sum_{i=1}^n \left\| A \left(\frac{\sum_{j=1}^n \exp(x_i^\top B x_j) x_j}{\sum_{j=1}^n \exp(x_i^\top B x_j)} \right) - y_i \right\|_2^2, \\
 R(\theta) &= \frac{1}{2} \|\Sigma^{1/2}(A^\top A - B^\top \Sigma B)\Sigma^{1/2}\|_F^2, \\
 R^{\text{emp}}(\theta) &= \frac{1}{2} \|\hat{\Sigma}^{1/2}(A^\top A - B^\top \hat{\Sigma} B)\hat{\Sigma}^{1/2}\|_F^2.
 \end{aligned}$$

We note that $L(\theta)$ and $L^{\text{emp}}(\theta)$ are both invariant under composition with $g(\theta, \lambda)$ for any λ . Recall that in the proof of Theorem 1 we define the residuals

$$\begin{aligned}
 \zeta_t^{\text{emp}} &= \nabla Q^{\text{emp}}(\bar{\theta}_t) - \nabla Q(\bar{\theta}_t), \\
 \hat{\zeta}_t &= \nabla \hat{Q}(\theta_t) - \nabla \tilde{Q}^{\text{emp}}(\theta_t).
 \end{aligned}$$

The proof of Theorem 1 requires that the variables are uniformly bounded in at all points of the gradient descent trajectory. Specifically, we prove:

Theorem 9 (Uniform concentration along gradient descent trajectory)

$$\sup_{t \in [m]} \max \left\{ \|\zeta_{t-1}^{\text{emp}}\|_F^2, \|\hat{\zeta}_{t-1}\|_F^2 \right\} \leq \frac{K \log^5(n) \log(mn/\delta)}{n}$$

for some constant $K \geq 0$.

The rest of this section is used to prove Theorem 9.

B.1. Proof of Theorem 9

We first show that

$$\sup_{t \in [m]} \|\zeta_{t-1}^{\text{emp}}\|_F^2 \leq \frac{K \log^5(n) \log(mn/\delta)}{n}$$

with probability at least $1 - \delta/2$, and then show that

$$\sup_{t \in [m]} \|\hat{\zeta}_{t-1}\|_F^2 \leq \frac{K \log^5(n) \log(mn/\delta)}{n}$$

with probability at least $1 - \delta/2$; the union bound implies the claim.

B.1.1. UNIFORM CONTROL OF $\|\zeta_t^{\text{emp}}\|_F^2$

Theorem 10 Fix any $\delta > 0$. Let $L^{\text{emp}}(\theta)$ be defined as in Lemma 14. With probability at least $1 - \delta$, the following inequalities holds for sufficiently large n :

$$\|\nabla L^{\text{emp}}(\theta) - \nabla L(\theta)\|_F^2 \leq \frac{K \log(1/\delta)}{n}$$

and

$$\|\nabla R^{\text{emp}}(\theta) - \nabla R(\theta)\|_F^2 \leq \frac{K \log(1/\delta)}{n}$$

for some $K \geq 0$.

Proof We observe that

$$\begin{aligned}\nabla L^{\text{emp}}(\theta) &= \begin{bmatrix} (A\hat{\Sigma}B^\top - M)\hat{\Sigma}B\hat{\Sigma} \\ \hat{\Sigma}(A\hat{\Sigma}B^\top - M)^\top A\hat{\Sigma} \end{bmatrix} \\ \nabla L(\theta) &= \begin{bmatrix} (A\Sigma B^\top - M)\Sigma B\Sigma \\ \Sigma(A\Sigma B^\top - M)^\top A\Sigma \end{bmatrix}.\end{aligned}$$

Let $\Delta = \hat{\Sigma} - \Sigma$. It is easy to see that

$$\|\nabla L^{\text{emp}}(\theta) - \nabla L(\theta)\|_F^2$$

can be upper bounded by a sum of terms, each of which has either linear, quadratic, or cubic dependence on $\|\Delta\|_F^2$. It is a standard result [8] that the empirical covariance exhibits exponential concentration around the true covariance, and in particular that $\|\Delta\|_F^2 \leq \frac{K \log(1/\delta)}{n}$ with probability $1 - \delta$. That implies the first part of the theorem. The second part of the theorem follows using an identical argument. ■

Corollary 11 For any $\delta \in (0, 1)$, the following inequality holds with probability at least $1 - \delta$:

$$\sup_{t \in [m]} \|\zeta_{t-1}^{\text{emp}}\|_F^2 \leq \frac{K \log(m/\delta)}{n}.$$

Proof Notice that

$$\|\zeta_t^{\text{emp}}\|_F^2 \leq 2\|\nabla L^{\text{emp}}(\bar{\theta}_t) - \nabla L(\bar{\theta}_t)\|_F^2 + 2\|\nabla R^{\text{emp}}(\bar{\theta}_t) - \nabla R(\bar{\theta}_t)\|_F^2.$$

Theorem 10 implies that each of these terms is bounded above by

$$\frac{K \log(m/\delta)}{n}$$

with probability $1 - \frac{\delta}{2m}$. Taking a union bound then implies the claim. ■

B.1.2. UNIFORM CONTROL OF $\|\hat{\zeta}_t\|_F^2$

Recall that

$$\hat{Q}(\theta) = \hat{L}(\theta) + \inf_{\lambda > 0} R^{\text{emp}}(g(\lambda, \theta))$$

and

$$\tilde{Q}^{\text{emp}}(\theta) = \inf_{\lambda > 0} (L^{\text{emp}}(g(\theta, \lambda)) + L^{\text{emp}}(g(\theta, \lambda))).$$

Notice that $L^{\text{emp}}(\theta)$ is invariant under the action of g , which implies that

$$Q^{\text{emp}}(\theta) = L^{\text{emp}}(\theta) + \inf_{\lambda > 0} R^{\text{emp}}(g(\lambda, \theta)).$$

It follows that

$$\hat{\zeta}_t = \nabla \hat{Q}(\theta_t) - Q^{\text{emp}}(\theta_t) = \nabla \hat{L}(\theta_t) - \nabla L^{\text{emp}}(\theta_t).$$

Also, recall that each of the $\{B_t\}_{t=1}^m$ is normalized to ensure that $\|B\|_F^2 \leq \gamma$.

We prove:

Theorem 12 *Let $\theta = (A, B)$ satisfy $\|B\|_F \leq \gamma$. With probability at least $1 - \delta$, the following inequality holds for all sufficiently large n :*

$$\|\nabla \hat{L}(\theta) - \nabla L^{\text{emp}}(\theta)\|_F^2 \leq \frac{K \log^5(n) \log(n/\delta)}{n}$$

where $K \geq 0$.

Proof Let $G_1(\theta)$ and $L^{\text{emp}}(\theta)$ be defined as in Lemma 13 and Lemma 14, respectively. Notice that

$$\nabla \hat{L}(\theta) - \nabla L(\theta) = (\nabla \hat{L}(\theta) - G_1(\theta)) + (G_1(\theta) - L^{\text{emp}}(\theta)) + (L^{\text{emp}}(\theta) - \nabla L(\theta)).$$

It follows that

$$\|\nabla \hat{L}(\theta_i) - \nabla L(\theta_i)\|_F^2 \leq 2\|\nabla \hat{L}(\theta_i) - G_1(\theta)\|_F^2 + 2\|G_1(\theta) - L^{\text{emp}}(\theta)\|_F^2 + 2\|L^{\text{emp}}(\theta) - \nabla L(\theta_i)\|_F^2.$$

In light of Lemmas (13–10), there exists some $K \geq 0$ such that each of these terms exceeds

$$\frac{K \log^{10}(n) \log(n/\delta)}{n}$$

with probability at most $\delta/3$, provided that n is sufficiently large. Applying the union bound yields the result. ■

Lemma 13 *Fix any $\delta > 0$ and set*

$$R = \max\left(1, \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n}\right), \quad p = \max(2, \log(3n/\delta)).$$

Define

$$G_1(\theta) = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n A \mathbb{E}_{-i}[\mu_i] \mathbb{E}_{-i}[\mu_i^\top] - \frac{1}{n} \sum_{i=1}^n M x_i \mathbb{E}_{-i}[\mu_i^\top] \\ \frac{1}{n} \sum_{i=1}^n x_i \mathbb{E}_{-i}[\mu_i^\top] A^\top A \mathbb{E}_{-i}[\Sigma_i] - \frac{1}{n} \sum_{i=1}^n x_i x_i^\top M^\top A \mathbb{E}_{-i}[\Sigma_i] \end{bmatrix}.$$

With probability at least $1 - \delta$, the following inequality holds:

$$\|\nabla \hat{L}(\theta) - G_1(\theta)\|_F^2 \leq \frac{KR^{10}p}{n}$$

for some $K \geq 0$ and n sufficiently large.

Proof Recall that

$$\nabla \hat{L}(\theta) = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n A \mu_i \mu_i^\top - \frac{1}{n} \sum_{i=1}^n M x_i \mu_i^\top \\ \frac{1}{n} \sum_{i=1}^n x_i \mu_i^\top A^\top A \Sigma_i - \frac{1}{n} \sum_{i=1}^n x_i x_i^\top M^\top A \Sigma_i \end{bmatrix}.$$

Applying Jensen's inequality, we observe that

$$\|\nabla \hat{L}(\theta) - G_1(\theta)\|_F^2 \leq \sum_{k=1}^4 D_k,$$

where we define

$$\begin{aligned}
 D_1 &= \frac{2}{n} \sum_{i=1}^n \left\| A\mu_i\mu_i^\top - A\mathbb{E}_{-i}[\mu_i]\mathbb{E}_{-i}[\mu_i^\top] \right\|_F^2, \\
 D_2 &= \frac{2}{n} \sum_{i=1}^n \left\| Mx_i\mu_i^\top - Mx_i\mathbb{E}_{-i}[\mu_i^\top] \right\|_F^2, \\
 D_3 &= \frac{2}{n} \sum_{i=1}^n \left\| x_i\mu_i^\top A^\top A\Sigma_i - x_i\mathbb{E}_{-i}[\mu_i^\top]A^\top A\mathbb{E}_{-i}[\Sigma_i] \right\|_F^2, \\
 D_4 &= \frac{2}{n} \sum_{i=1}^n \left\| x_i x_i^\top M^\top A\Sigma_i - x_i x_i^\top M^\top A\mathbb{E}_{-i}[\Sigma_i] \right\|_F^2.
 \end{aligned}$$

We bound each of these terms separately. We first bound D_1 . We observe that for each $i \in [n]$,

$$\mu_i\mu_i^\top - \mathbb{E}_{-i}[\mu_i]\mathbb{E}_{-i}[\mu_i^\top] = (\mu_i - \mathbb{E}_{-i}[\mu_i])(\mu_i - \mathbb{E}_{-i}[\mu_i])^\top + (\mu_i - \mathbb{E}_{-i}[\mu_i])\mathbb{E}_{-1}[\mu_i^\top] + \mathbb{E}_{-1}[\mu_i](\mu_i - \mathbb{E}_{-1}[\mu_i])^\top.$$

It follows that

$$\begin{aligned}
 D_1 &\leq \frac{6}{n} \|A\|_F^2 \sum_{i=1}^n \|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_F^4 + \frac{12}{n} \|A\|_F^2 \sum_{i=1}^n \|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_F^2 \|\mathbb{E}_{-i}[\mu_i]\|_F^2 \\
 &\leq \frac{6}{n} \|A\|_F^2 \sum_{i=1}^n \|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_F^4 + \frac{12}{n} \|A\|_F^2 \sum_{i=1}^n \|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_F^2 \mathbb{E}_{-i}[M_n^2],
 \end{aligned}$$

where we applied Jensen's inequality and the fact that $\|\mu_i\|_2 \leq M_n$ because μ_i is a convex combination of the covariates. We now bound D_2 . Applying submultiplicativity of the Frobenius norm, it is immediate that

$$D_2 \leq \frac{2}{n} \|M\|_F^2 \sum_{i=1}^n \|x_i\|_2^2 \|\mu_i^\top - \mathbb{E}_{-i}[\mu_i^\top]\|_2^2.$$

We now bound D_3 . We observe that

$$\mu_i^\top A^\top A\Sigma_i - \mathbb{E}_{-i}[\mu_i^\top]A^\top A\mathbb{E}_{-i}[\Sigma_i]$$

can be rewritten as

$$(\mu_i - \mathbb{E}_{-i}[\mu_i])^\top A^\top A(\Sigma_i - \mathbb{E}_{-i}[\Sigma_i]) + (\mu_i - \mathbb{E}_{-i}[\mu_i])^\top A^\top A\mathbb{E}_{-i}[\Sigma_i] + \mathbb{E}_{-i}[\mu_i^\top]A^\top A(\Sigma_i - \mathbb{E}_{-i}[\Sigma_i]).$$

It follows that

$$\begin{aligned}
 D_3 &\leq \frac{3}{n} \|A\|_F^4 \sum_{i=1}^n \|x_i\|_2^2 \left[\|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_2^2 \|\Sigma_i - \mathbb{E}_{-i}[\Sigma_i]\|_F^2 + \|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_2^2 \|\mathbb{E}_{-i}[\Sigma_i]\|_F^2 + \|\mathbb{E}_{-i}[\mu_i]\|_2^2 \|\Sigma_i - \mathbb{E}_{-i}[\Sigma_i]\|_F^2 \right] \\
 &\leq \frac{3}{n} \|A\|_F^4 \sum_{i=1}^n \|x_i\|_2^2 \left[\|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_2^2 \|\Sigma_i - \mathbb{E}_{-i}[\Sigma_i]\|_F^2 + \|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_2^2 \mathbb{E}_{-i}[M_n^4] + \mathbb{E}_{-i}[M_n^2] \|\Sigma_i - \mathbb{E}_{-i}[\Sigma_i]\|_F^2 \right]
 \end{aligned}$$

where we once again applied Jensen's inequality and the fact that Σ_i is a convex combination of the rank-1 outer products $\{x_i x_i^\top\}_{i=1}^n$. Lastly, we bound D_4 . Applying submultiplicativity of the Frobenius norm, it is immediate that

$$D_4 \leq \frac{2}{n} \|M\|_F^2 \|A\|_F^2 \sum_{i=1}^n \|x_i\|_2^4 \|\Sigma_i^\top - \mathbb{E}_{-i}[\Sigma_i^\top]\|_F^2.$$

Let R, K be defined as in Theorem 22, and set $p = \max(2, \log(3n/\delta))$. Let \mathcal{A}_i denote the event that $\|x_i\|_2 \leq R$, let \mathcal{B}_i denote the event that

$$\|\mu_i - \mathbb{E}_{-i}[\mu_i]\|_2^2 \leq \frac{KR^4 p}{n},$$

and let \mathcal{C}_i denote the event that

$$\|\Sigma_i - \mathbb{E}_{-i}[\Sigma_i]\|_F^2 \leq \frac{KR^6 p}{n}.$$

Notice that Lemma 24 implies that on the event $\mathcal{D} = \bigcap_{i=1}^n (A_i \cap B_i \cap C_i)$, each of the terms $\{D_k\}_{k=1}^4$ is bounded above by $KR^{10}pn^{-1}$ for some constant K provided that n is sufficiently large. We now show that \mathcal{D} fails to occur with probability at most δ . Notice that

$$\mathcal{D}^c = \bigcup_{i=1}^n \mathcal{A}_i^c \cup (\mathcal{B}_i^c \cap \mathcal{A}_i) \cup (\mathcal{C}_i^c \cap \mathcal{A}_i).$$

Applying the union bound, we see that

$$\begin{aligned} \Pr(\mathcal{D}^c) &\leq \sum_{i=1}^n [\Pr(\mathcal{A}_i^c) + \Pr(\mathcal{B}_i^c \mid \mathcal{A}_i) \Pr(\mathcal{A}_i) + \Pr(\mathcal{C}_i^c \mid \mathcal{A}_i) \Pr(\mathcal{A}_i)] \\ &\leq \sum_{i=1}^n [\Pr(\mathcal{A}_i^c) + \Pr(\mathcal{B}_i^c \mid \mathcal{A}_i) + \Pr(\mathcal{C}_i^c \mid \mathcal{A}_i)]. \end{aligned}$$

In light of Theorems 22 and 23 and Lemma 25, we see that this sum is at most $2\delta/3 + n^{-3}$, and in particular is bounded above by δ provided that $n \geq (3/\delta)^3$. \blacksquare

Lemma 14 Fix any $\delta > 0$ and set

$$R = \max\left(1, \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n}\right), \quad p = \max(2, \log(3n/\delta)).$$

Let $G_1(\theta)$ be defined as in Lemma 13. Recall that

$$\nabla L^{\text{emp}}(\theta) = \begin{bmatrix} (A\Sigma B^\top - M) \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top\right) B\Sigma \\ \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top\right) (A\Sigma B^\top - M)^\top A\Sigma \end{bmatrix}.$$

With probability at least $1 - \delta$, the following inequality holds:

$$\|G_1(\theta) - L^{\text{emp}}(\theta)\|_F^2 \leq \frac{KR^{10}p}{n}$$

for some $K \geq 0$ and n sufficiently large.

Proof Consulting the definitions of $G_1(\theta)$ and $L^{\text{emp}}(\theta)$ and applying Jensen's inequality, we observe that

$$\|G_1(\theta) - L^{\text{emp}}(\theta)\|_F^2 \leq \sum_{k=1}^4 E_k,$$

where we define

$$\begin{aligned}
 E_1 &= \frac{2}{n} \|A\|_F^2 \sum_{i=1}^n \left\| \mathbb{E}_{-i}[\mu_i] \mathbb{E}_{-i}[\mu_i^\top] - \Sigma B^\top x_i x_i^\top B \Sigma \right\|_F^2, \\
 E_2 &= \frac{2}{n} \|M\|_F^2 \sum_{i=1}^n \|x_i\|_2^2 \left\| \mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right\|_2^2, \\
 E_3 &= \frac{2}{n} \sum_{i=1}^n \|x_i\|_2^2 \left\| \mathbb{E}_{-i}[\mu_i^\top] A^\top A \mathbb{E}_{-i}[\Sigma_i] - x_i^\top B \Sigma A^\top A \Sigma \right\|_2^2, \\
 E_4 &= \frac{2}{n} \|M\|_F^2 \|A\|_F^2 \sum_{i=1}^n \|x_i\|_2^4 \left\| \mathbb{E}_{-i}[\Sigma_i] - \Sigma \right\|_F^2.
 \end{aligned}$$

We bound each of these terms separately. We first bound E_1 . We observe that for each $i \in [n]$,

$$\mathbb{E}_{-i}[\mu_i] \mathbb{E}_{-i}[\mu_i^\top] - \Sigma B^\top x_i x_i^\top B \Sigma$$

can be rewritten as

$$\left(\mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right) \left(\mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right)^\top + \left(\mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right) x_i^\top B \Sigma + \Sigma B^\top x_i \left(\mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right)^\top.$$

It follows that

$$E_1 \leq \frac{6}{n} \|A\|_F^2 \sum_{i=1}^n \left\| \mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right\|_2^4 + \frac{12}{n} \|A\|_F^2 \|B\|_F^2 \|\Sigma\|_F^2 \sum_{i=1}^n \left\| \mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right\|_2^2 \|x_i\|_2^2.$$

The terms E_2 and E_4 are already in a convenient form, so we now bound E_3 . We observe that for each $i \in [n]$,

$$\mathbb{E}_{-i}[\mu_i^\top] A^\top A \mathbb{E}_{-i}[\Sigma_i] - x_i^\top B \Sigma A^\top A \Sigma$$

can be rewritten as

$$\left(\mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right)^\top A^\top A \left(\mathbb{E}_{-i}[\Sigma_i] - \Sigma \right) + \left(\mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right)^\top A^\top A \Sigma + x_i^\top B \Sigma A^\top A \left(\mathbb{E}_{-i}[\Sigma_i] - \Sigma \right)^\top.$$

It follows that E_3 is at most the sum of

$$\frac{6}{n} \|A\|_F^4 \sum_{i=1}^n \left\| \mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right\|_2^2 \left\| \mathbb{E}_{-i}[\Sigma_i] - \Sigma \right\|_F^2$$

and

$$\frac{6}{n} \|A\|_F^4 \sum_{i=1}^n \left\| \mathbb{E}_{-i}[\mu_i] - \Sigma B^\top x_i \right\|_2^2 \left(\|\Sigma\|_F^2 + \left\| \mathbb{E}_{-i}[\Sigma_i] - \Sigma \right\|_F^2 \right) \|B\|_F^2 \|\Sigma\|_F^2 \|x_i\|_2^2.$$

Let \mathcal{A}_i denote the event that $\|x_i\|_2 \leq R$. Notice that Lemma 18 and Lemma 19 together imply that on the event $\cap_{i=1}^n \mathcal{A}_i$, each of the terms $\{E_k\}_{k=1}^4$ is deterministically bounded above by $K R^{10} p n^{-1}$ for some constant K provided that n is sufficiently large. We now show that the event $\cap_{i=1}^n \mathcal{A}_i$ fails to occur with probability at most δ . Applying the union bound, we see that the failure probability is at most $\sum_{i=1}^n \Pr(\mathcal{A}_i^c)$. In light of Lemma 25, we see that this sum is at most n^{-3} , and in particular is bounded above by δ provided that $n \geq (1/\delta)^3$. \blacksquare

B.2. Conditional moments and conditional variance of softmax weights

We use the following classical inequality to establish concentration of the softmax-weighted mean and softmax-weighted covariance.

Theorem 15 (Pisier's inequality) *Let f be a smooth, symmetric function of x_1, \dots, x_n , where each x_i is sampled i.i.d. from $\mathcal{N}(0, \Sigma)$. The following moment bound holds:*

$$\mathbb{E} \|f(x_1, \dots, x_n) - \mathbb{E}[f(x_1, \dots, x_n)]\|_F^p \leq C_p (n-1)^{p/2} \|\Sigma\|_{\text{op}}^{p/2} \mathbb{E} \|\nabla_{x_1} f(x_1, \dots, x_n)\|_F^p$$

Lemma 16 (Some useful gradients) *The following identities hold:*

$$\begin{aligned} \nabla_{x_2} \pi_{12} &= \pi_{12} (1 - \pi_{12}) B^\top x_1, \\ \nabla_{x_2} \pi_{11} &= -\pi_{11} \pi_{12} B^\top x_1, \\ \nabla_{x_2} \mu_1 &= \pi_{12} (x_2 - \mu_1) x_1^\top B + \pi_{12} I, \\ \nabla_{x_2} \Sigma_1 &= \pi_{12} \left(\left((x_2 - \mu_1)(x_2 - \mu_1)^\top - \Sigma_1 \right) \otimes (B^\top x_1) \right) + \pi_{12} \left(I \otimes (x_2 - \mu_1) + (x_2 - \mu_1) \otimes I \right). \end{aligned}$$

Proof These identities are easily verified via direct calculation. ■

Lemma 17 (Conditional moments of softmax weights) *Fix any $p \geq 1$. The following inequalities hold:*

$$\begin{aligned} \mathbb{E}_{-1}[\pi_{12}^p] &\leq \frac{1}{(n-1)^p} \exp\left(\frac{p^2}{2} x_1^\top B \Sigma B^\top x_1\right) \\ \mathbb{E}_{-1}[\pi_{11}^p] &\leq \frac{1}{(n-1)^p} \exp\left(p x_1^\top B x_1\right) \exp\left(\frac{p^2}{2(n-1)} x_1^\top B \Sigma B^\top x_1\right). \end{aligned}$$

Proof To prove the first inequality, we observe that

$$\begin{aligned} \pi_{12} &= \frac{\exp(x_1^\top B x_2)}{\sum_{j=1}^n \exp(x_1^\top B x_j)} \\ &\leq \frac{\exp(x_1^\top B x_2)}{\sum_{j=2}^n \exp(x_1^\top B x_j)} \\ &\leq \frac{1}{n-1} \exp\left(\frac{n-2}{n-1} x_1^\top B x_2\right) \exp\left(-\frac{1}{n-1} x_1^\top B \sum_{j=3}^n x_j\right), \end{aligned}$$

where we applied the AM-GM inequality in the last step. It follows that

$$\begin{aligned} \mathbb{E}_{-1}[\pi_{12}^p] &= \frac{1}{(n-1)^p} \mathbb{E}_{-1} \left[\exp\left(\frac{p(n-2)}{n-1} x_1^\top B x_2\right) \exp\left(-\frac{p}{n-1} x_1^\top B \sum_{j=3}^n x_j\right) \right] \\ &= \frac{1}{(n-1)^p} \mathbb{E}_2 \left[\exp\left(\frac{p(n-2)}{n-1} x_1^\top B x_2\right) \right] \mathbb{E}_{-1,-2} \left[\exp\left(-\frac{p}{n-1} x_1^\top B \sum_{j=3}^n x_j\right) \right], \end{aligned}$$

where we used the fact that x_2 is independent of $\{x_j\}_{j=3}^n$. We observe that

$$\mathbb{E}_2 \left[\exp \left(\frac{p(n-2)}{n-1} x_1^\top B x_2 \right) \right] = \exp \left(\frac{p^2(n-2)^2}{2(n-1)^2} x_1^\top B \Sigma B^\top x_1 \right)$$

and

$$\mathbb{E}_{-1,-2} \left[\exp \left(-\frac{p}{n-1} x_1^\top B \sum_{j=3}^n x_j \right) \right] = \exp \left(\frac{p^2(n-2)}{2(n-1)^2} x_1^\top B \Sigma B^\top x_1 \right).$$

Putting the pieces together, we see that

$$\mathbb{E}_{-1}[\pi_{12}^p] \leq \frac{1}{(n-1)^p} \exp \left(\frac{p^2}{2} x_1^\top B \Sigma B^\top x_1 \right),$$

where we used the numerical fact that

$$\frac{(n-2)^2}{(n-1)^2} + \frac{n-2}{(n-1)^2} \leq 1.$$

We now prove the second inequality using a similar calculation. We observe that

$$\begin{aligned} \pi_{11} &= \frac{\exp(x_1^\top B x_1)}{\sum_{j=1}^n \exp(x_1^\top B x_j)} \\ &\leq \frac{\exp(x_1^\top B x_1)}{\sum_{j=2}^n \exp(x_1^\top B x_j)} \\ &\leq \frac{1}{n-1} \exp(x_1^\top B x_1) \exp \left(-\frac{1}{n-1} x_1^\top B \sum_{j=2}^n x_j \right), \end{aligned}$$

where we applied the AM-GM inequality in the last step. It follows that

$$\mathbb{E}_{-1}[\pi_{11}^p] = \frac{1}{(n-1)^p} \exp(p x_1^\top B x_1) \mathbb{E}_{-1} \left[\exp \left(-\frac{p}{n-1} x_1^\top B \sum_{j=2}^n x_j \right) \right].$$

We observe that

$$\mathbb{E}_{-1} \left[\exp \left(-\frac{p}{n-1} x_1^\top B \sum_{j=2}^n x_j \right) \right] = \exp \left(\frac{p^2}{2(n-1)} x_1^\top B \Sigma B^\top x_1 \right),$$

which completes the proof. ■

Lemma 18 *Set*

$$R = \max \left(1, \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n} \right).$$

Fix $x_1 \in \mathbb{R}^d$ such that $\|x_1\|_2 \leq R$. The following inequality holds:

$$\|\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1\|_2^2 \leq \frac{KR^2}{n^2},$$

where $K \geq 0$ is a constant, for sufficiently large n .

Proof We see that

$$\begin{aligned} \left\| \mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1 \right\|_F^2 &= \left\| \mathbb{E}_{-1} \left[\sum_{j=1}^n \pi_{1j} x_j \right] - \Sigma B^\top x_1 \right\|_2^2 \\ &= \left\| \mathbb{E}_{-1} \left[\sum_{j=1}^n \pi_{1j} (x_j - \Sigma B^\top x_1) \right] \right\|_2^2, \end{aligned}$$

where we used the fact that the softmax weights sum to unity. Using the fact that the covariates are i.i.d., we see that this expression is at most

$$2 \left\| \mathbb{E}_{-1}[\pi_{11}] (I - \Sigma B^\top) x_1 \right\|_2^2 + 2 \left\| (n-1) \mathbb{E}_{-1} \left[\pi_{12} (x_2 - \Sigma B^\top x_1) \right] \right\|_2^2. \quad (14)$$

We bound each of the two terms of (14) separately. Applying homogeneity of the Frobenius norm and Lemma 17, we see that the first term is at most

$$\frac{2}{(n-1)^2} \exp \left(\frac{2}{n-1} x_1^\top B \Sigma B^\top x_1 \right) \exp(2x_1^\top B x_1) \|I - \Sigma B^\top\|_{\text{op}}^2 \|x_1\|_2^2.$$

Applying the assumptions that $\|x_1\|_2 \leq R$ and the assumption that $\|B\|_{\text{op}} \leq \gamma$, we see that this term is bounded above by KR^2n^{-2} . We now bound the second term of (14). Applying Gaussian integration by parts, we observe that

$$\mathbb{E}_{-1} \left[\pi_{12} (x_2 - \Sigma B^\top x_1) \right] = -\mathbb{E}_{-1}[\pi_{12}^2] \Sigma B^\top x_1.$$

It follows that the second term of (14) is at most

$$2(n-1)^2 (\mathbb{E}_{-1}[\pi_{12}^2])^2 \left\| \Sigma B^\top x_1 \right\|_F^2.$$

Applying Lemma 17, we see that this is at most

$$\frac{2}{(n-1)^2} \exp(4x_1^\top B \Sigma B^\top x_1) \|\Sigma B^\top\|_{\text{op}}^2 \|x_1\|_2^2.$$

Applying our assumptions on x_1 and B once more, we see that this term is at most KR^2n^{-2} . ■

Lemma 19 *Set*

$$R = \max \left(1, \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n} \right).$$

Fix $x_1 \in \mathbb{R}^d$ such that $\|x_1\|_2 \leq R$. The following inequality holds:

$$\left\| \mathbb{E}_{-1}[\Sigma_1] - \Sigma \right\|_F^2 \leq \frac{KR^8}{n^2},$$

where $K \geq 0$ is a constant.

Proof Recall that

$$\Sigma_1 = \sum_{j=1}^n \pi_{1j} x_j x_j^\top - \mu_1 \mu_1^\top.$$

It follows that

$$\|\mathbb{E}_{-1}[\Sigma_1] - \Sigma\|_F^2 \leq 2 \left\| \mathbb{E}_{-1} \left[\sum_{j=1}^n \pi_{1j} x_j x_j^\top \right] - (\Sigma + \Sigma B^\top x_1 x_1^\top B \Sigma) \right\|_F^2 + 2 \left\| \mathbb{E}_{-1}[\mu_1 \mu_1^\top] - \Sigma B^\top x_1 x_1^\top B \Sigma \right\|_F^2. \quad (15)$$

We bound each of these two terms separately.

We first bound the first term of (15). Using the fact that the covariates are i.i.d., we see that

$$\mathbb{E}_{-1} \left[\sum_{j=1}^n \pi_{1j} x_j x_j^\top \right] = \mathbb{E}_{-1}[\pi_{11}] x_1 x_1^\top + (n-1) \mathbb{E}_{-1}[\pi_{12} x_2 x_2^\top].$$

By Gaussian integration by parts,

$$\mathbb{E}_{-1}[\pi_{12} x_2 x_2^\top] = \mathbb{E}_{-1}[\pi_{12}] \Sigma + \mathbb{E}_{-1}[\pi_{12}(1-\pi_{12})(1-2\pi_{12})] \Sigma B^\top x_1 x_1^\top B \Sigma, \quad (16)$$

where we used the identities

$$\nabla_{x_2} \pi_{12} = \pi_{12}(1-\pi_{12}) B^\top x_1, \quad \nabla_{x_2}^2 \pi_{12} = \pi_{12}(1-\pi_{12})(1-2\pi_{12}) B^\top x_1 x_1^\top B.$$

The softmax weights sum to one, which implies the elementary identity

$$\mathbb{E}_{-1}[\pi_{11}] + (n-1) \mathbb{E}_{-1}[\pi_{12}] = 1. \quad (17)$$

Applying (16) and (17), we see that the expression

$$\mathbb{E}_{-1} \left[\sum_{j=1}^n \pi_{1j} x_j x_j^\top \right] - (\Sigma + \Sigma B^\top x_1 x_1^\top B \Sigma)$$

can be rewritten as

$$\mathbb{E}_{-1}[\pi_{11}](x_1 x_1^\top - \Sigma) + ((n-1) \mathbb{E}_{-1}[\pi_{12}(1-\pi_{12})(1-2\pi_{12})] - 1) \Sigma B^\top x_1 x_1^\top B \Sigma.$$

We observe that

$$\pi_{12}(1-\pi_{12})(1-2\pi_{12}) = \pi_{12} - 3\pi_{12}^2 + 2\pi_{12}^3,$$

which implies that

$$(n-1) \mathbb{E}_{-1}[\pi_{12}(1-\pi_{12})(1-2\pi_{12})] - 1 = -\mathbb{E}_{-1}[\pi_{11}] - 3(n-1) \mathbb{E}_{-1}[\pi_{12}^2] + 2(n-1) \mathbb{E}_{-1}[\pi_{12}^3],$$

where we applied the identity (17). Applying Lemma 17 and the scaling assumptions on B , we see that for sufficiently large n , the following inequalities hold:

$$\mathbb{E}_{-1}[\pi_{11}] \leq \frac{K}{n}, \quad \mathbb{E}_{-1}[\pi_{12}^2] \leq \frac{K}{n^2}, \quad \mathbb{E}_{-1}[\pi_{12}^3] \leq \frac{K}{n^3}.$$

It follows that

$$|(n-1) \mathbb{E}_{-1}[\pi_{12}(1-\pi_{12})(1-2\pi_{12})] - 1| \leq \frac{K}{n}.$$

Putting the pieces together, we see that

$$\left\| \mathbb{E}_{-1} \left[\sum_{j=1}^n \pi_{1j} x_j x_j^\top \right] - \Sigma - \Sigma B^\top x_1 x_1^\top B \Sigma \right\|_F^2 \leq \frac{K}{n^2} \left(\|x_1 x_1^\top - \Sigma\|_F^2 + \|\Sigma B^\top x_1 x_1^\top B \Sigma\|_F^2 \right).$$

Using the assumption that $\|x_1\|_2 \leq R$ and the fact that $R \geq 1$, we obtain the bound

$$\left\| \mathbb{E}_{-1} \left[\sum_{j=1}^n \pi_{1j} x_j x_j^\top \right] - \Sigma - \Sigma B^\top x_1 x_1^\top B \Sigma \right\|_F^2 \leq \frac{KR^4}{n^2}.$$

We now bound the second term of (15). We observe that

$$\mathbb{E}_{-1}[\mu_1 \mu_1^\top] - \Sigma B^\top x_1 x_1^\top B \Sigma$$

can be rewritten as the sum of two terms, namely

$$\mathbb{E}_{-1} \left[(\mu_1 - \mathbb{E}_{-1}[\mu_1])(\mu_1 - \mathbb{E}_{-1}[\mu_1])^\top \right] + (\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1)(\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1)^\top$$

and

$$(\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1)(\Sigma B^\top x_1)^\top + \Sigma B^\top x_1(\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1)^\top.$$

Applying Jensen's inequality, we see that

$$\left\| \mathbb{E}_{-1}[\mu_1 \mu_1^\top] - \Sigma B^\top x_1 x_1^\top B \Sigma \right\|_F^2 \leq 4\mathbb{E}_{-1} \left[\|\mu_1 - \mathbb{E}_{-1}[\mu_1]\|_2^4 \right] + 8\|\Sigma B^\top x_1\|_2^2 \|\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1\|_2^2 + 4\|\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1\|_2^2$$

Applying Lemma 20 with $p = 4$, along with Lemma 24 and the fact that $R \geq 1$, we see that

$$\mathbb{E}_{-1} \left[\|\mu_1 - \mathbb{E}_{-1}[\mu_1]\|_2^4 \right] \leq \frac{KR^8}{n^2}.$$

Lemma 18 immediately implies that

$$\|\mathbb{E}_{-1}[\mu_1] - \Sigma B^\top x_1\|_2^2 \leq \frac{KR^2}{n^2}.$$

Furthermore, it is clear that

$$\|\Sigma B^\top x_1\|_2^2 \leq KR^2.$$

These inequalities together imply the bound

$$\left\| \mathbb{E}_{-1}[\mu_1 \mu_1^\top] - \Sigma B^\top x_1 x_1^\top B \Sigma \right\|_F^2 \leq \frac{KR^8}{n^2}$$

for sufficiently large n . ■

Appendix C. Concentration of the softmax mean and softmax covariance

Lemma 20 *The following inequality holds for all $p \geq 2$:*

$$\mathbb{E}_{-1} \left\| \mu_1 - \mathbb{E}_{-1}[\mu_1] \right\|_F^p \leq C_p \tau_1^{2p^2} 2^{p-1} (n-1)^{-p/2} \|\Sigma\|_{\text{op}}^{p/2} \left(2^p \|B\|_{\text{op}}^p (\mathbb{E}_{-1}[M_n^{4p}])^{1/2} + d^{p/2} \right),$$

where $\tau_1 = \exp\left(\frac{1}{2} x_1^\top B \Sigma B^\top x_1\right)$.

Proof

Pisier's inequality yields the bound

$$\mathbb{E}_{-1} \|\mu_1 - \mathbb{E}_{-1}[\mu_1]\|_F^p \leq C_p (n-1)^{p/2} \|\Sigma\|_{\text{op}}^{p/2} \mathbb{E}_{-1} \|\nabla_{x_2} \mu_1\|_F^p, \quad (18)$$

where we used the fact that the covariates $\{x_j\}_{j=2}^n$ are i.i.d. We recall from Lemma 16 that

$$\nabla_{x_2} \mu_1 = \pi_{12} (x_2 - \mu_1) x_1^\top B + \pi_{12} I.$$

We see that

$$\|\nabla_{x_2} \mu_1\|_F^p \leq 2^{p-1} \pi_{12}^p \left(\|x_2 - \mu_1\|_F^p \|B\|_{\text{op}}^p \|x_1\|_2^p + d^{p/2} \right). \quad (19)$$

Set $M_n = \sup_{i \in [n]} \|x_i\|_2$. It is clear that $\|x_1\|_2 \leq M_n$. Applying the triangle inequality and using the fact that μ_1 is a convex combination of the covariates, we see that $\|x_2 - \mu_1\|_2 \leq 2M_n$. Plugging these bounds into (19), we obtain the bound

$$\|\nabla_{x_2} \mu_1\|_F^p \leq 2^{p-1} \pi_{12}^p \left(2^p \|B\|_{\text{op}}^p M_n^{2p} + d^{p/2} \right). \quad (20)$$

Plugging this bound into (18) and applying the Cauchy-Schwarz inequality, we obtain the bound

$$\mathbb{E}_{-1} \|\mu_1 - \mathbb{E}_{-1}[\mu_1]\|_F^p \leq C_p 2^{p-1} (n-1)^{p/2} \|\Sigma\|_{\text{op}}^{p/2} \left(\mathbb{E}_{-1}[\pi_{12}^{2p}] \right)^{1/2} \left(2^p \|B\|_{\text{op}}^p \left(\mathbb{E}_{-1}[M_n^{4p}] \right)^{1/2} + d^{p/2} \right).$$

Applying Lemma 17, we obtain the claim. \blacksquare

Lemma 21 *The following inequality holds for all $p \geq 2$:*

$$\mathbb{E}_{-1} \|\Sigma_1 - \mathbb{E}_{-1}[\Sigma_1]\|_F^p \leq C_p \tau_1^{2p^2} 4^{p-1} (n-1)^{-p/2} \|\Sigma\|_{\text{op}}^{p/2} \left(2^{2p+1} \|B\|_{\text{op}}^p \left(\mathbb{E}_{-1} M_n^{6p} \right)^{1/2} + 2^{p+1} d^{p/2} \left(\mathbb{E}_{-1} M_n^{2p} \right)^{1/2} \right),$$

where $\tau_1 = \exp\left(\frac{1}{2} x_1^\top B \Sigma B^\top x_1\right)$.

Proof Pisier's inequality yields the bound

$$\mathbb{E}_{-1} \|\Sigma_1 - \mathbb{E}_{-1}[\Sigma_1]\|_F^p \leq C_p (n-1)^{p/2} \|\Sigma\|_{\text{op}}^{p/2} \mathbb{E}_{-1} \|\nabla_{x_2} \Sigma_1\|_F^p, \quad (21)$$

where we used the fact that the covariates $\{x_j\}_{j=2}^n$ are i.i.d. We recall from Lemma 16 that

$$\nabla_{x_2} \Sigma_1 = \pi_{12} \left(\left((x_2 - \mu_1)(x_2 - \mu_1)^\top - \Sigma_1 \right) \otimes (B^\top x_1) \right) + \pi_{12} \left(I \otimes (x_2 - \mu_1) + (x_2 - \mu_1) \otimes I \right).$$

We see that

$$\|\nabla_{x_2} \Sigma_1\|_F^p \leq 4^{p-1} \pi_{12}^p \left(\|B\|_{\text{op}}^p \|x_1\|_2^p \|x_2 - \mu_1\|_F^{2p} + \|B\|_{\text{op}}^p \|x_1\|_2^p \|\Sigma_1\|_F^p + 2d^{p/2} \|x_2 - \mu_1\|_2^p \right). \quad (22)$$

Set $M_n = \sup_{i \in [n]} \|x_i\|_2$. It is clear that $\|x_1\|_2 \leq M_n$. Applying the triangle inequality and using the fact that μ_1 is a convex combination of the covariates, we see that $\|x_2 - \mu_1\|_2 \leq 2M_n$. Similarly, convexity of the function $\|\cdot\|_F^p$ implies that

$$\|\Sigma_1\|_F^p \leq \sup_{i \in [n]} \left\| (x_i - \mu_1)(x_i - \mu_1)^\top \right\|_F^p$$

which in turn implies the bound $\|\Sigma_1\|_F^p \leq 4^p M_n^{2p}$. Plugging these bounds into (22), we obtain the bound

$$\|\nabla_{x_2} \Sigma_1\|_F^p \leq C_p 4^{p-1} \pi_{12}^p \left(2^{2p+1} \|B\|_{\text{op}}^p M_n^{3p} + 2^{p+1} d^{p/2} M_n^p \right).$$

Plugging this bound into (21) and applying the Cauchy-Schwarz inequality, we obtain the bound

$$\mathbb{E}_{-1} \left\| \Sigma_1 - \mathbb{E}_{-1}[\Sigma_1] \right\|_F^p \leq C_p 4^{p-1} (n-1)^{p/2} \|\Sigma\|_{\text{op}}^{p/2} \left(\mathbb{E}_{-1} \pi_{12}^{2p} \right)^{1/2} \left(2^{2p+1} \|B\|_{\text{op}}^p \left(\mathbb{E}_{-1} M_n^{6p} \right)^{1/2} + 2^{p+1} d^{p/2} \left(\mathbb{E}_{-1} M_n^{2p} \right)^{1/2} \right).$$

Applying Lemma 17, we obtain the claim. \blacksquare

Theorem 22 *Set*

$$R = \max \left(1, \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n} \right), \quad p = \max(2, \log(1/\delta)).$$

Let γ be defined as in Lemma 26 and assume that $\|B\|_F \leq \gamma$. Fix $\delta \in (0, 1)$ and $x_1 \in \mathbb{R}^d$ such that $\|x_1\|_2 \leq R$. There exist constants K and N with polynomial dependence on M, Σ such that

$$\Pr \left(\left\| \mu_1 - \mathbb{E}_{-1}[\mu_1] \right\|_2^2 \geq \frac{KR^4 p}{n} \right) \leq \delta,$$

provided that $n \geq N$.

Proof Define

$$Z = \left\| \mu_1 - \mathbb{E}_{-1}[\mu_1] \right\|_2.$$

Applying Markov's inequality, we observe that for any $t \geq 0$ and any $p \geq 1$,

$$\begin{aligned} \Pr(Z^2 \geq t) &= \Pr(Z \geq \sqrt{t}) \\ &\leq \left(\frac{\|Z\|_p}{\sqrt{t}} \right)^p, \end{aligned}$$

where we define

$$\|Z\|_p = (\mathbb{E}_{-1}[Z^p])^{1/p}.$$

Set $p = \max(2, \log(1/\delta))$. Suppose that $\|Z\|_p \leq \frac{1}{3}\sqrt{K}R^2\sqrt{\frac{p}{n}}$ when n is sufficiently large, for some $K \geq 0$. Setting $t = KR^4\frac{p}{n}$, we see that

$$\begin{aligned} \Pr(Z \geq \sqrt{t}) &\leq 3^{-p} \\ &\leq 3^{\log(\delta)} \\ &= \delta^{\log 3} \\ &\leq \delta, \end{aligned}$$

where we used the fact that $\delta < 1$ and $\log 3 > 1$. We now show that when n is sufficiently large, $\|Z\|_p \leq \frac{1}{3}\sqrt{K}R^2\sqrt{\frac{p}{n}}$ for some $K \geq 0$. Applying Lemma 20, we see that

$$\|Z\|_p^p \leq C_p \tau_1^{2p^2} 2^{p-1} (n-1)^{-p/2} \|\Sigma\|_{\text{op}}^{p/2} \left(2^p \|B\|_{\text{op}}^p \left(\mathbb{E}_{-1}[M_n^{4p}] \right)^{1/2} + d^{p/2} \right),$$

where $\tau_1 = \exp\left(\frac{1}{2}x_1^\top B \Sigma B^\top x_1\right)$. It is clear that under the assumption that $\|x_1\|_2 \leq R$,

$$\begin{aligned} \tau_1^{2p} &\leq \exp(pR^2\|B\|_{\text{op}}^2\|\Sigma\|_{\text{op}}) \\ &\leq e, \end{aligned}$$

where we used the assumption that $\|B\|_{\text{op}} \leq \gamma$.

Recall from Lemma 24 that

$$\mathbb{E}_{-1}[M_n^k] \leq \left(R + \sqrt{\text{Tr}(\Sigma)} + \sqrt{\|\Sigma\|_{\text{op}}(\sqrt{2\log n} + \sqrt{k+2})} \right)^k.$$

Assume that n is chosen such that $\sqrt{2\log n} \geq \sqrt{4p+2}$. We see that

$$\mathbb{E}_{-1}[M_n^{4p}] \leq (2R)^{4p}.$$

We also recall the elementary estimates $C_p^{1/p} \leq \frac{\pi}{2}\sqrt{p}$ and $(n-1)^{-1/2} \leq 2n^{-1}$ for $n \geq 2$. Putting the pieces together and using the fact that $R \geq 1$, we see that

$$\begin{aligned} \|Z\|_p &\leq 2\pi e \|\Sigma\|_{\text{op}}^{1/2} \max(2\|B\|_{\text{op}}, \sqrt{d}) R^2 \sqrt{\frac{p}{n}} \\ &= \frac{1}{3} \sqrt{K} R^2 \sqrt{\frac{p}{n}}, \end{aligned}$$

where we set

$$K = 9 \cdot 2^2 \pi^2 e^2 \|\Sigma\|_{\text{op}} \max(4\|B\|_{\text{op}}^2, d).$$

■

Theorem 23 *Set*

$$R = \max\left(1, \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n}\right), \quad p = \max(2, \log(1/\delta)).$$

Let γ be defined as in Lemma 26 and assume that $\|B\|_F \leq \gamma$. Fix $\delta \in (0, 1)$ and $x_1 \in \mathbb{R}^d$ such that $\|x_1\|_2 \leq R$. There exist constants K and N with polynomial dependence on M, Σ such that

$$\Pr\left(\|\Sigma_1 - \mathbb{E}_{-1}[\Sigma_1]\|_F^2 \geq \frac{KR^6 p}{n}\right) \leq \delta,$$

provided that $n \geq N$.

Proof Define

$$Z = \|\Sigma_1 - \mathbb{E}_{-1}[\Sigma_1]\|_F.$$

Applying Markov's inequality, we observe that for any $t \geq 0$ and any $p \geq 1$,

$$\begin{aligned} \Pr(Z^2 \geq t) &= \Pr(Z \geq \sqrt{t}) \\ &\leq \left(\frac{\|Z\|_p}{\sqrt{t}}\right)^p, \end{aligned}$$

where we define

$$\|Z\|_p = (\mathbb{E}_{-1}[Z^p])^{1/p}.$$

Set $p = \max(2, \log(1/\delta))$. Suppose that $\|Z\|_p \leq \frac{1}{3}\sqrt{K}R^3\sqrt{\frac{p}{n}}$ when n is sufficiently large, for some $K \geq 0$. Setting $t = KR^6\frac{p}{n}$, we see that

$$\begin{aligned} \Pr\left(Z \geq \sqrt{t}\right) &\leq 3^{-p} \\ &\leq 3^{\log(\delta)} \\ &= \delta^{\log 3} \\ &\leq \delta, \end{aligned}$$

where we used the fact that $\delta < 1$ and $\log 3 > 1$. We now show that when n is sufficiently large, $\|Z\|_p \leq \frac{1}{3}\sqrt{K}R^3\sqrt{\frac{p}{n}}$ for some $K \geq 0$. Applying Lemma 21, we see that

$$\|Z\|_p^p \leq C_p \tau_1^{2p^2} 4^{p-1} (n-1)^{-p/2} \|\Sigma\|_{\text{op}}^{p/2} \left(2^{2p+1} \|B\|_{\text{op}}^p (\mathbb{E}_{-1}[M_n^{6p}])^{1/2} + 2^{p+1} d^{p/2} (\mathbb{E}_{-1}[M_n^{2p}])^{1/2} \right),$$

where $\tau_1 = \exp\left(\frac{1}{2}x_1^\top B \Sigma B^\top x_1\right)$. It is clear that under the assumption that $\|x_1\|_2 \leq R$,

$$\begin{aligned} \tau_1^{2p} &\leq \exp(pR^2 \|B\|_{\text{op}}^2 \|\Sigma\|_{\text{op}}) \\ &\leq e, \end{aligned}$$

where we used the assumption that $\|B\|_{\text{op}} \leq \gamma$.

Recall from Lemma 24 that

$$\mathbb{E}_{-1}[M_n^k] \leq \left(R + \sqrt{\text{Tr}(\Sigma)} + \sqrt{\|\Sigma\|_{\text{op}}(\sqrt{2\log n} + \sqrt{k+2})} \right)^k.$$

Notice that this bound increases monotonically in k . Assume that n is chosen such that $\sqrt{2\log n} \geq \sqrt{6p+2}$. We see that

$$\max(\mathbb{E}_{-1}[M_n^{2p}], \mathbb{E}_{-1}[M_n^{6p}]) \leq (2R)^{6p}.$$

We also recall the elementary estimates $C_p^{1/p} \leq \frac{\pi}{2}\sqrt{p}$ and $(n-1)^{-1/2} \leq 2n^{-1}$ for $n \geq 2$. Putting the pieces together and using the fact that $R \geq 1$, we see that

$$\begin{aligned} \|Z\|_p &\leq 32\pi e \|\Sigma\|_{\text{op}}^{1/2} \max(2\|B\|_{\text{op}}, \sqrt{d}) R^3 \sqrt{\frac{p}{n}} \\ &= \frac{1}{3}\sqrt{K}R^3\sqrt{\frac{p}{n}}, \end{aligned}$$

where we set

$$K = 9 \cdot 32^2 \pi^2 e^2 \|\Sigma\|_{\text{op}} \max(4\|B\|_{\text{op}}^2, d).$$

■

Lemma 24 *Let $R \geq 0$ and fix $x_1 \in \mathbb{R}^d$ such that $\|x_1\|_2 \leq R$. Let $M_n = \sup_{i \in [n]} \|x_i\|_2$. The following inequality holds:*

$$\mathbb{E}_{-1}[M_n^k] \leq \left(R + \sqrt{\text{Tr}(\Sigma)} + \sqrt{\|\Sigma\|_{\text{op}}(\sqrt{2\log n} + \sqrt{k+2})} \right)^k.$$

In particular, if n is chosen such that $2\log n \geq k+2$, and R is chosen such that

$$R \geq \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n},$$

then the following inequality holds:

$$\mathbb{E}_{-1}[M_n^k] \leq (2R)^k.$$

Proof It is clear that

$$M_n \leq R + \max_{2 \leq i \leq n} \|x_i\|_2.$$

Applying Minkowski's inequality, we see that

$$\left(\mathbb{E}_{-1}[M_n^k]\right)^{1/k} \leq R + \left(\mathbb{E}\left[\max_{2 \leq i \leq n} \|x_i\|_2^k\right]\right)^{1/k}.$$

Gaussian concentration of Lipschitz functions [13] lets us easily bound the second term:

$$\left(\mathbb{E}\left[\max_{2 \leq i \leq n} \|x_i\|_2^k\right]\right)^{1/k} \leq \mathbb{E}\left[\max_{2 \leq i \leq n} \|x_i\|_2\right] + \sqrt{\|\Sigma\|_{\text{op}}(k+2)}.$$

The Gaussian maximal inequality yields the bound

$$\mathbb{E}\left[\max_{2 \leq i \leq n} \|x_i\|_2\right] \leq \sqrt{\text{Tr}(\Sigma)} + \sqrt{2\|\Sigma\|_{\text{op}} \log n}.$$

The claim follows immediately. ■

Lemma 25 *Let R be chosen such that*

$$R \geq \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n},$$

and let $x \sim \mathcal{N}(0, \Sigma)$. The following tail bound holds:

$$\Pr(\|x\|_2 \geq R) \leq n^{-4}.$$

Proof This claim follows immediately from the standard Gaussian tail bound

$$\Pr(\|x\|_2 \geq \sqrt{\text{Tr}(\Sigma)} + t) \leq \exp\left(-\frac{t^2}{2\|\Sigma\|_{\text{op}}}\right).$$
■

Appendix D. Initialization

Lemma 26 *Fix any $\delta \in (0, 1)$. Set*

$$\varepsilon_\Sigma = \min\left\{\|\Sigma\|_{\text{op}}^{-1}, \frac{\alpha}{8\sqrt{\beta}}\|\Sigma\|_{\text{op}}^{-1/2}\right\}, \quad \varepsilon_\lambda = \lambda_{\max},$$

where λ_{\max} is defined as in Theorem 7. Let $\alpha, \beta, \varepsilon_0$, and S be defined as in Theorem ??, and let \bar{S} denote the set of points which are ε_0 -close to S in P -norm. The probability that any of the following events fail to occur is at most δ , provided that n is sufficiently large:

(E_1) *The empirical covariance is close to the true covariance, and the inverse of the empirical covariance is close to the inverse of the true covariance:*

$$\max\left\{\|\hat{\Sigma} - \Sigma\|_{\text{op}}, \|\hat{\Sigma}^{-1} - \Sigma^{-1}\|_{\text{op}}\right\} \leq \varepsilon_\Sigma.$$

(E₂) The initial choice of parameters θ_0 lies in $\bar{\mathcal{S}}$.

(E₃) The random function $\hat{\lambda}^*$ is uniformly close to λ^* , and $\hat{\lambda}^{-*}$ is uniformly close to λ^{-*} :

$$\sup_{(A,B) \in \bar{\mathcal{S}}} \max \left\{ |\hat{\lambda}^*(A, B) - \lambda^*(A, B)|, |\hat{\lambda}^{-*}(A, B) - \lambda^{-*}(A, B)| \right\} \leq \varepsilon_\lambda.$$

(E₄) The renormalization factor $\hat{\gamma}$ is at most γ , where we define

$$\begin{aligned} \nu &= \max \left(1, \sqrt{\text{Tr}(\Sigma)} + 2\sqrt{2\|\Sigma\|_{\text{op}} \log n} \right), \\ \hat{\nu} &= \max \left(1, \sqrt{\text{Tr}(\hat{\Sigma})} + 2\sqrt{2\|\hat{\Sigma}\|_{\text{op}} \log n} \right), \\ p &= \max(2, \log(1/\delta)) \end{aligned}$$

and set

$$\gamma = \sqrt{\frac{\|\Sigma\|_{\text{op}}}{p\nu^2}}, \quad \hat{\gamma} = \frac{1}{2} \sqrt{\frac{\|\hat{\Sigma}\|_{\text{op}}}{p\hat{\nu}^2}}.$$

Proof We show that each event E_i occurs with probability tending to 1 as n increases; the union bound implies that the probability that any event fails to occur tends to zero, and in particular can be made at most δ by taking n sufficiently large.

(E₁) It is a standard result that

$$\max \left\{ \|\hat{\Sigma} - \Sigma\|_F^2, \|\hat{\Sigma}^{-1} - \Sigma^{-1}\|_F^2 \right\} \leq K \max \{ \|\Sigma\|_{\text{op}}^2, \|\Sigma^{-1}\|_{\text{op}}^2 \} \frac{d + \log(1/\delta)}{n}$$

with probability at least $1 - \delta$, for some $K \geq 0$ and sufficiently large n . We refer to [8, 13] for details.

(E₂) Recall that the parameters are initialized at the point

$$\theta_0 = (A_0, B_0) = \begin{bmatrix} \hat{U}\hat{\Gamma}^{1/2}\hat{\Sigma}^{-1/2} \\ \hat{\Sigma}^{-1/2}\hat{V}\hat{\Gamma}^{1/2}\hat{\Sigma}^{-1/2} \end{bmatrix},$$

where $\hat{U}\hat{\Gamma}\hat{V}^\top$ is a singular value decomposition of $\hat{M}\hat{\Sigma}^{1/2}$ and

$$\hat{M} = \frac{1}{n} \sum_{i=1}^n y_i x_i^\top \hat{\Sigma}^{-1}.$$

Using the fact that $y_i = Mx_i + z_i$ for each $i \in [n]$, we see that

$$\begin{aligned} \hat{M} &= \frac{1}{n} \sum_{i=1}^n \left(Mx_i x_i^\top + z_i x_i^\top \right) \hat{\Sigma}^{-1} \\ &= \frac{1}{n} \sum_{i=1}^n \left(Mx_i x_i^\top + z_i x_i^\top \right) \hat{\Sigma}^{-1} \\ &= M + \frac{1}{n} \sum_{i=1}^n z_i x_i^\top \hat{\Sigma}^{-1}. \end{aligned}$$

It follows that

$$\|M - \hat{M}\|_{\text{op}} \leq \left\| \frac{1}{n} \sum_{i=1}^n z_i x_i^\top \hat{\Sigma}^{-1} \right\|_{\text{op}}.$$

It is clear that \hat{M} converges to M in probability, because the random variables $\{z_i\}_{i=1}^n$ are i.i.d. zero mean subgaussian random variables which are independent of the covariates, so

$$\left\| \frac{1}{n} \sum_{i=1}^n z_i x_i^\top \hat{\Sigma}^{-1} \right\|_{\text{op}}$$

converges to zero in probability. It follows that $\hat{M}\hat{\Sigma}^{1/2}$ converges to $M\Sigma^{1/2}$ in probability. Lemma 5.14 from [11] implies that the factors $\hat{U}, \hat{\Gamma}, \hat{V}$ converge (up to orthogonal transformations) to U, Γ, V , where $U\Gamma V^\top$ is a singular value decomposition of $M\Sigma^{1/2}$. It follows that $\|\theta_0 - \theta_0^*\|_P^2$ converges in probability to zero, where θ_0^* is the projection of θ_0 onto \mathcal{S} .

(E₃) Recall that λ^* is defined in terms of the true covariance Σ :

$$\lambda^*(A, B) = \left(\frac{\|\Sigma^{1/2} A^\top A \Sigma^{1/2}\|_F}{\|\Sigma^{1/2} B^\top B \Sigma^{1/2}\|_F} \right)^{1/4}.$$

Similarly, $\hat{\lambda}^*$ is defined in terms of the empirical covariance $\hat{\Sigma}$:

$$\hat{\lambda}^*(A, B) = \left(\frac{\|\hat{\Sigma}^{1/2} A^\top A \hat{\Sigma}^{1/2}\|_F}{\|\hat{\Sigma}^{1/2} B^\top B \hat{\Sigma}^{1/2}\|_F} \right)^{1/4}.$$

In Lemma 7, we show that $\lambda^*(\theta)$ is continuous on $\bar{\mathcal{S}}$. Note that $\bar{\mathcal{S}}$ is compact. Applying the continuous mapping theorem, we see that $\hat{\lambda}^*(\theta)$ converges uniformly to $\lambda^*(\theta)$ on $\bar{\mathcal{S}}$.

(E₄) An identical argument shows that $\hat{\gamma}$ converges in probability to $\frac{1}{2}\gamma$ as n increases, and in particular satisfies $\hat{\gamma} \leq \gamma$ with probability $1 - \delta$ for sufficiently large n .

■