

High-Dimensional Statistical Inference from Coarse and
Nonlinear Data: Algorithms and Guarantees

Dissertation

Presented in Partial Fulfillment of the Requirements for the Degree
Doctor of Philosophy in the Graduate School of The Ohio State
University

By

Haoyu Fu, M.S.

Graduate Program in Electrical and Computer Engineering

The Ohio State University

2019

Dissertation Committee:

Yingbin Liang, Advisor

Yuejie Chi, Co-Advisor

Kiryung Lee

Philip Schniter

© Copyright by

Haoyu Fu

2019

Abstract

Learning a postulated parametric model from the acquired data to extract useful information is of great importance in modern signal processing, machine learning and statistics. The linear model, where the observed data are assumed to depend linearly on the input data, has been studied extensively and applied successfully to many applications. However, the linear assumption is quite restricted, creating a major roadblock for its accuracy and universality, since the dependency of the data is nonlinear in general, and cannot be approximated by a linear model. The challenges of learning these nonlinear models include high computational cost and susceptibility to local minima in their associated optimization problems. Through case studies, we highlight the statistical and computational issues when learning from high-dimensional coarse and nonlinear data, with the hope of shedding light on resolving these challenges.

In this thesis, we consider data from typical signal processing and machine learning applications. In the context of signal processing, we study the problem of estimating spectrally-sparse signals from their quantized noisy complex-valued random linear measurements, a problem arising naturally from analog-to-digital conversion in sub-Nyquist spectrum sampling. We first study the effects of quantization on estimating the spectrum by characterizing the Cramér-Rao bound under the additive white Gaussian noise. We use the calculated bound to highlight the trade-off between the sample complexity and the bit depth under different signal-to-noise ratios for a

fixed budget of bits. Secondly, we formulate a convex optimization approach based on atomic norm soft thresholding to estimate the spectrum of the signal, which is computationally more efficient than the maximum-likelihood estimator.

Moving to the context of machine learning, we study several one-hidden-layer neural network models for nonlinear regression using both cross-entropy and least-squares loss functions. The neural-network-based models have attracted a significant amount of research interest due to the success of deep learning in practical domains such as computer vision and natural language processing. Learning such neural-network-based models often requires solving a non-convex optimization problem. We propose different strategies to characterize the optimization landscape of the non-convex loss functions and provide guarantees on the statistical and computational efficiency of optimizing these loss functions via gradient descent.

Dedicated to my parents

Acknowledgments

First and foremost, I would like to express my sincere gratitudes to my two advisors, Prof. Yingbin Liang and Prof. Yuejie Chi, for their great guidance and support over the past few years. Prof. Liang and Prof. Chi are both great researchers. Their passion for exploring unknown problems has motivated and inspired me a lot while I develop my own research path. Their valuable feedback and insightful suggestions help me a lot in improving the quality of my research work. Their kindness and wisdom would continue to impact me. I feel so lucky to be their student.

I would like to thank my thesis committee members, Prof. Kiryung Lee and Prof. Philip Schniter, for their insightful comments and invaluable suggestions on my thesis.

I would like to thank Dr. Pu (Perry) Wang who was my mentor when I interned at Mitsubishi Electric Research Laboratories. He gave me plenty guidance and meaningful suggestions on our research project.

I'm also grateful to fellow graduate students in The Ohio State University, including Yuanxin Li, Liming Wang, Ziwei Guan, Yi Zhou and Vince Monardo for their generous help, encouragement and interesting discussions. I want to thank the faculty and staff from The Ohio State University for their selfless dedication to the education career. I benefit tremendously from their solid professionalism.

Lastly, I would like to thank my parents, and my sister for their love, support and encouragement. I thank Yubo for sharing hundreds of delightful stories, and bringing so much happiness to my life. To them I dedicate this thesis.

Vita

- 2014B.Eng., The measurement and control
technology and instrument,
Wuhan University of Technology,
Wuhan, China
- 2019M.S., Electrical and Computer Engi-
neering,
The Ohio State University,
Columbus, USA

Publications

Research Publications

H. Fu, Y. Chi, and Y. Liang, “Guaranteed Recovery of One-Hidden-Layer Neural Networks via Cross Entropy”, *submitted*.

H. Fu and Y. Chi, “Quantized Spectral Compressed Sensing: CramerRao Bounds and Recovery Algorithms”, *IEEE Transactions on Signal Processing*, vol. 66, no. 12, pp. 3268-3279, 2018.

Y. Chi and H. Fu, “Subspace Learning From Bits”, *IEEE Transactions on Signal Processing*, vol. 65, no. 12, pp. 4429-4442, 2017.

H. Fu, Y. Chi, and Y. Liang, “Local Geometry of Cross Entropy Loss in Learning One-Hidden-Layer Neural Networks”, *IEEE International Symposium on Information Theory*, Paris, France, 2019.

H. Fu, P. Wang, T. Koike-Akino, R. Ma, B.Wang, P.V. Orlik, W. Tsujitaz, K. Sadamotoz, Y. Sawaz, K. Katox and M. Nakajima, “Terahertz Imaging of Multi-Level Pseudo-Random Reflectance”, *International Conference on Infrared, Millimeter, and Terahertz Waves*, Nagoya, Japan, 2018.

P. Wang, H. Fu, T. Koike-Akino, and P.V. Orlik, “Multi-Layer Terahertz Imaging of Non-Overlapping Contents”, *IEEE Sensor Array and Multi-Channel Signal Processing Workshop*, Sheffield, UK, 2018.

H. Fu, P. Wang, T. Koike-Akino, P.V. Orlik, and Y. Chi, “Terahertz imaging of binary reflectance with variational Bayesian inference”, *IEEE International Conference on Acoustics, Speech and Signal Processing*, Alberta, Canada, 2018.

H. Fu and Y. Chi, “Compressive Spectrum Estimation using Quantized Measurements”, *Asilomar Conference on Signals, Systems, and Computers*, Pacific Grove, USA, 2017.

H. Fu and Y. Chi, “Principal Subspace Estimation for Low-rank Toeplitz Covariance Matrices with Binary Sensing”, *Asilomar Conference on Signals, Systems, and Computers*, Pacific Grove, USA, 2016.

Fields of Study

Major Field: Electrical and Computer Engineering

Table of Contents

	Page
Abstract	ii
Dedication	iv
Acknowledgments	v
Vita	vii
List of Figures	xii
1. Introduction	1
1.1 Motivation	1
1.2 Convex and Nonconvex Approaches	3
1.3 Spectrum Estimation from Quantized Measurements	5
1.4 Learning One-Hidden-Layer Neural Networks for Binary Classification	9
1.5 Guaranteed Recovery of CNN with ReLU Activations	15
1.6 Notations	18
2. Line Spectrum Estimation with Quantized Measurements	19
2.1 Problem Formulation	20
2.2 Cramer-Rao Bounds and Trade-offs	22
2.2.1 CRB for 1-Bit Quantization	23
2.2.2 CRB for General Quantization	24
2.2.3 Numerical Evaluations of CRB	28
2.3 Atomic Norm Soft Thresholding for Quantized Spectral Compressed Sensing	29
2.3.1 Backgrounds on Atomic Norms	29
2.3.2 Atomic Soft-Thresholding with Quantized Measurements	30

2.3.3	Performance Guarantees	31
2.4	Extension to the Multiple Vector Case	33
2.5	Numerical Experiments	35
2.5.1	Single Vector Case	35
2.5.2	Multiple Vector Case	40
3.	Learning One-Hidden-Layer Neural Network for Binary Classification	42
3.1	Problem Formulation	43
3.2	Gradient Descent and its Performance Guarantee	45
3.2.1	Uniform local strong convexity	45
3.2.2	Performance Guarantees of GD	48
3.3	Initialization via Tensor Method	50
3.3.1	Preliminary and Algorithm	50
3.3.2	Performance Guarantee of Initialization	53
3.4	Numerical Experiments	54
4.	Guaranteed Recovery of CNN with ReLU Activations	58
4.1	Problem Formulation	58
4.2	Minimizing the Population Risk	60
4.3	Minimizing the Empirical Risk	64
4.4	Numerical Experiments	67
5.	Conclusion and Future Work	69
5.1	Concluding Remarks	69
5.2	Future Work	71
	Appendices	74
A.	Proofs for Chapter 2	74
A.1	Proof of Theorem 2	74
A.2	Proof of Lemma 3	76
B.	Proofs for Chapter 3	78
B.1	Gradient and Hessian of the Population Loss	78
B.1.1	The FCN case	78
B.1.2	The CNN case	79
B.2	Proof of Theorem 3	79

B.3	Proof of Theorem 4	82
B.4	Proof of Auxiliary Lemmas	85
B.4.1	Proof of Lemma 5.	85
B.4.2	Proof of Lemma 6	90
B.4.3	Proof of Lemma 7	95
B.4.4	Proof of Lemma 8	100
B.4.5	Proof of Lemma 9	101
B.4.6	Proof of Lemma 11	102
B.4.7	Proof of Lemma 13	104
B.4.8	Proof of Lemma 12	104
B.4.9	Proof of Lemma 14	105
B.4.10	Proof of Lemma 15	106
B.5	Proof of Theorem 5	107
C.	Proofs for Chapter 4	110
C.1	Preliminary	110
C.2	Proof of GD on the Population Risk	113
C.2.1	Proof of Theorem 6	113
C.3	Proof of GD on the Empirical Risk	118
C.3.1	Proof of Lemma 1	118
C.3.2	Proof of Auxiliary Lemmas	121
	Bibliography	135

List of Figures

Figure	Page
2.1 The CRB under different bit-depths with respect to SNR for a fixed number of measurements $m = 100$. Here, $n = 64$ and $K = 3$. Each row represents the CRB for estimating the frequency, amplitude and phase of one spectral atom.	26
2.2 The CRB under different bit-depths with respect to SNR for a fixed number of bits $B = 100$. In this case, 2-bit quantization only has half the number of measurements of the 1-bit case. Here, $n = 64$ and $K = 3$. Each row represents the CRB for estimating the frequency, amplitude and phase of one spectral atom.	27
2.3 The value of λ with respect to SNR before quantization.	32
2.4 Frequency localization via peaks of the dual polynomial, superimposed on the ground truth.	36
2.5 Normalized reconstruction error with respect to the number of measurements at different SNRs with or without quantization.	37
2.6 Normalized reconstruction error with respect to the spectral sparsity level at different SNRs before quantization.	38
2.7 Mean square error of frequency localization with respect to SNR using 1-bit measurements, CRB is provided as a benchmark: (a) first frequency; (b) second frequency.	39
2.8 Performance with respect to the number of snapshots at different SNRs using 1-bit measurements: (a) signal reconstruction error; (b) frequency estimation error measured in Hausdorff distance.	40

3.1	Illustration of two types of one-hidden-layer neural networks considered in this Chapter: (a) a fully-connected network (FCN); (b) a non-overlapping convolutional neural network (CNN).	44
3.2	Illustration of $\rho(\sigma)$ for both FCN and CNN with the sigmoid activation. 49	
3.3	For FCN (3.1) fix $K = 3$. (a) Success rate of converging to the same local minima with respect to the sample complexity for various d with threshold 10^{-4} ; (b) Average estimation error of gradient descent in a local neighborhood of the ground truth with respect to the sample complexity for various d . The x-axis is scaled to illuminate the correct scaling between n and d	55
3.4	For CNN (3.2), fix $K = 3$. (a) Success rate of converging to the same local minima with respect to the sample complexity for various d with threshold 10^{-14} ; (b) Average estimation error of gradient descent in a local neighborhood of the ground truth with respect to the sample complexity for various d . The x-axis is scaled to illuminate the correct scaling between n and d	56
4.1	Illustration of a one-hidden-layer convolutional neural network without overlap	59
4.2	Contour plot of the population risk function (4.5) when $m = 2$ and $K = 3$	61
4.3	For CNN, fix $K = 10$, $m = 35$ and $d = 350$, the NMSE with respect to the number of steps of gradient descent for various n . (a) local initialization; (b) Gaussian initialization.	68
C.1	Numerical integral with respect to ρ	133

Chapter 1: Introduction

This thesis studies several problems on learning from coarse and nonlinear data. In this chapter, we first explain the motivation of these problems, and describe the major approach to solving them. We then summarize our main contributions with respect to the state-of-the-art in the literature.

1.1 Motivation

In order to extract latent information from the collected data samples in many machine learning, signal processing and statistical inference tasks, a parametric model is generally proposed to postulate a data generating mechanism. Fitting an effective and reasonable model is then equivalent to recovering a finite set of parameters. In this thesis, we follow this direction and focus on parametric models to fit the data. Assuming the observed data is generated from a true parametric model, our goal is to recover the parameters of the underlying model from the acquired data.

One example is the classical linear model, where we aim to estimate a vector $\mathbf{w}^* \in \mathbb{R}^d$ from a set of independent and identically distributed (i.i.d.) data samples

$$\{(\mathbf{x}_i, y_i) : \mathbf{x}_i \in \mathbb{R}^d \text{ and } y_i \in \mathbb{R}\}_{i=1}^n,$$

by assuming the observation y_i depends linearly on the vector \mathbf{w}^* , through the following model:

$$y_i = \mathbf{x}_i^\top \mathbf{w}^* + \nu_i, \quad i = 1, 2, \dots, n, \quad (1.1)$$

where ν_i stands for some noise.

To find an estimator $\hat{\mathbf{w}}$ that approximates \mathbf{w}^* , a common strategy is to form a proper optimization problem with the data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, i.e., solving the following optimization problem

$$\begin{aligned} \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \quad L_n(\mathbf{w}) &= \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}_i, y_i), \\ \text{subject to} \quad \mathbf{w} &\in \mathcal{C}. \end{aligned} \quad (1.2)$$

Here, the loss function $L_n(\mathbf{w})$ measures the discrepancy between the candidate model and the true model, and the set \mathcal{C} is the feasible region which has incorporated the constraint or prior information, if any. For example, when the noise ν_i follows i.i.d. Gaussian distribution and $\mathcal{C} = \mathbb{R}^d$, it is not a bad idea to take $\ell(\mathbf{w}; \mathbf{x}_i, y_i) = (\mathbf{x}_i^\top \mathbf{w} - y_i)^2$, the least-squares, or the quadratic loss, which yields the maximum likelihood estimator (MLE).

Although the linear model (1.1) is equipped with many good properties, for example, it is easy to interpret, and computationally efficient, to name a few. However, the stringent assumption that the output depends on the input via a linear relationship has been a serious impediment for its accuracy and universality, since the relationship of the data can be nonlinear in general. One example of the non-linearity is when the observed data is quantized into a binary observation, i.e., $y_i \in \{-1, 1\}$. Such binary observations have emerged in many applications, such as 1-bit compressed sensing [1] in signal processing, and binary classification in machine learning. Another example

is that the observed data is continuous but does not necessarily depend linearly on the input data, such as the output of a neural network with nonlinear activations.

The binary or the general nonlinear nature of the output data has raised many new challenges which are not inherent in the linear model (1.1). In particular, the computational cost of the MLE with quantized data could be much higher. When the quantizer is not given, the effect of quantization has not been fully understood yet, a MLE cannot be formulated. Finding the MLE for a general nonlinear model often requires solving a non-convex optimization problem which is much more challenging. Hence, in this thesis we study the problem of learning from coarse and nonlinear data to shed light on understanding quantization effects as well as developing algorithms for solving a type of non-convex optimization problem.

1.2 Convex and Nonconvex Approaches

One of the most popular methods to recover the model parameterized with \mathbf{w}^* from the observed data is via solving the optimization problem (1.2). The loss function $L_n(\cdot)$ or the feasible region \mathcal{C} is often nonconvex in many problems. The main challenge caused by the nonconvexity is the existence of local minima and saddle points that can attract the algorithm.

One way to circumvent this challenge is to reformulate the problem as a convex optimization problem via convex relaxation, which excuses the need to worry about the spurious local minima. For example, in the problem of spectrum estimation which will be studied later, we assume the signal \mathbf{w}^* is approximately sparse in the spectral domain. However, incorporating such sparse prior makes the feasible set \mathcal{C} nonconvex, and searching in such a nonconvex space is not computationally tractable. To deal

with this problem, we take an atomic norm [2] constraint which convexifies the exact sparsity constraint on \mathbf{w} . Such a convex relaxation approach has achieved tremendous theoretical success in solving various problems [3, 4].

Differently from convex relaxation, an alternative approach is to directly solve the nonconvex optimization problem with careful local landscape analysis and initialization. Although solving general nonconvex optimization problems is computationally intractable due to the existence of spurious local minima or saddle points, solving certain nonconvex problems under proper statistical models may not be that difficult. For certain types of nonconvex problems, it was shown that in the local region around the global minima, the objective function usually possesses a benign landscape [5–7]. For example, in the problem of recovering neural-network-based models that will be studied later, we found that there often exists a reasonably large basin of attraction around the ground truth \mathbf{w}^* . Particularly, we can either explore the second-order (Hessian) property of the loss function to establish the local strong convexity or check the first-order (gradient) property to verify a so-called regularity condition under some mild conditions. Once initialized in such local regions, simple algorithms like gradient descent can find the global optimum \mathbf{w}^* .

This Thesis explores the convex approach to solve the problem of spectrum estimation from quantized measurements in the context of signal processing, and adopts the nonconvex approach to solve the neural-network-based model recovery problem in the context of machine learning.

1.3 Spectrum Estimation from Quantized Measurements

We first study a problem in the context of signal processing, i.e., estimating a spectrally-sparse signals from its heavy quantized linear measurements. In this problem, the non-linearity of the data comes from the quantization procedure, i.e., the output data depends on the input through

$$y_i = f(\mathbf{x}_i^\top \mathbf{w}^* + \nu_i), \quad (1.3)$$

where $f(\cdot)$ is used to model the quantization non-linearity. More specifically, we study high-resolution spectrum estimation of a band-limited signal from quantization of its noisy random linear measurements. The signals of interest are spectrally sparse, which are modeled as a linear superposition of complex sinusoids with continuous-valued frequencies. In the extreme 1-bit case, the quantization is based on the quadrants of the complex-valued measurements. More generally, sophisticated quantization schemes such as Lloyds quantizer [8] can be used to allow a higher bit depth. The specific form of the quantizer can be either known or unknown. In addition, the quantized measurements may be additionally contaminated by a noise model, in order to model imperfections in the quantization.

The reasons we study this problem are two-fold. Compressed Sensing (CS) [9, 10] has emerged as an effective approach to allow sub-Nyquist sampling [11–13] when the wideband signal is approximately sparse in the spectral domain. The resulting paradigm is referred to as *Compressive Spectrum Sensing* [14, 15]. Significant focus has been put on reducing the sampling rates of the analog-to-digital converters (ADC), which only covers one aspect of the operations of ADCs. Quantization, which maps the analog samples into a finite number of bits for digital processing, is

another necessary step that requires careful treatments. Most existing works, with a few exceptions, assume that the samples are quantized at a high bit level so that the quantization error is relatively small and well-behaved. Another motivation comes from the application wideband spectrum sensing in bandwidth-constrained wireless networks [16, 17]. In order to reduce the communication overhead, each sensor transmits quantized messages, e.g. 1-bit messages; and it is necessary to estimate wideband spectrum from quantized measurements at the fusion center. Moreover, the quantization scheme might be unknown, due to lack of the knowledge of noise statistics or privacy constraints. Therefore, it is necessary to develop estimators that do not require exact knowledge of the quantizers.

Hence the study aims at understanding the fundamental limits of quantization, as well as developing computationally efficient algorithms, for compressive spectrum sensing and parameter estimation, in particular in the regime of *heavy* quantization where it is no longer appropriate to model quantization errors as bounded additive noise. Examining the figure-of-merit of ADCs, two key specifications are the sampling rate and the effective number of bits (ENOB), which is the number of bits per measurement, also known as the *bit depth*. Typically, a small bit depth allows a high sampling rate, and vice versa [18]. Therefore, it is critical to understand the fundamental trade-off between sampling rate and bit depth for high-resolution spectrum estimation. Though the importance of understanding such trade-off has been realized in the context of CS [19, 20], they haven't been studied for the task of *parameter estimation* using estimation-theoretic tools.

Contributions

In this study, we first derive the Cramér-Rao bound (CRB) for estimating multiple frequencies and their complex amplitudes assuming additive white Gaussian noise (AWGN) and the Lloyd’s quantizer, using a fixed and deterministic CS measurement matrix. Our bounds suggest that the CRB experiences a phase transition depending on the signal-to-noise ratio (SNR) *before* quantization. In the low SNR regime it is *noise-limited*, and behaves similarly as if there was no quantization; in the high SNR regime, it is *quantization-limited*, and experiences severe performance degeneration due to quantization. Furthermore, we use the derived CRB to answer the following question: given the same budget of bits, should we use more measurements (high sample complexity) with low bit-depth, or fewer measurements (low sample complexity) with high bit-depth? We answer this question by comparing 1-bit versus 2-bit quantization schemes using the CRB, and demonstrate the answer depends on the SNR. At low SNR, 1-bit measurements are preferred, while at high SNR, 2-bit measurements are preferred.

It is well-known that maximum likelihood estimators approach the performance of CRB asymptotically at high SNR [21], however, their implementation requires exact knowledge of the likelihood function, which in our problem, includes the exact form of the quantizer and noise statistics. However, such knowledge may not be available in certain applications. Therefore, our goal is to develop estimators that do not require the knowledge of the quantization scheme. To mitigate basis mismatch [22], atomic norm [2, 23–30] has been proposed recently to promote spectral sparsity via convex optimization without discretizing the frequencies onto a finite grid, which has found

applications in signal denoising, interpolation of missing data, and frequency localization of spectrally-sparse signals. Existing atomic norm minimization algorithms assume unquantized measurements that are possibly contaminated by additive noise, and a direct application will lead to highly sub-optimal performance when a significant amount of the quantized measurements *saturate* [31].

We propose a novel atomic norm soft thresholding (AST) algorithm [29] to recover spectrally-sparse signals and estimate the frequencies from their 1-bit quantized measurements. Our algorithm is based on finding the proximal mapping of properly designed surrogate signals, that are formed by linear combinations of the sample-modulated measurement vectors, with respect to the atomic norm to promote spectral sparsity. In other words, we aim to find signals that balance between the proximity to the surrogate signals and the small atomic norm. Moreover, the frequencies can be localized without knowing the model order a priori, by examining the peak of a dual polynomial constructed from the dual solution. Alternatively, conventional subspace methods can be used to estimate the frequencies using the recovered spectral signal. The proposed algorithm can be generalized to handle quantizations of noisy random linear measurements of multiple spectrally-sparse signals [25], where each signal contains the same set of frequencies with different coefficients. The proposed algorithms do not require knowledge of the specific form of the quantizer, and therefore can be applied even when the quantizer is unknown.

Related Work

This study is closely related to 1-bit compressed sensing [1, 32–38], which aims to recover a sparse signal from signs of random linear measurements. In particular, Plan and Vershynin [34–37] generalize this idea to reconstructing signals that belong

to some low-dimensional set. Very recently, [39] studied a similar setup and proposed a new algorithm using projected gradient descent. The surrogate signals used in our algorithm can be traced back to [20, 37]. The difference lies in that instead of projecting the surrogate signals directly onto some low-dimensional set, we adopt the proximal mapping of the surrogate signals with respect to the atomic norm. Several algorithms have been proposed in the CS literature to deal with general quantization schemes [19, 40] and nonlinear measurement schemes [37, 39], however the focus has been on reconstruction of sparse signals in a finite dictionary, whereas our focus is on parameter estimation and reconstructing sparse signals in a parametric dictionary containing an infinite number of atoms.

There are also several conflicting evidence regarding the trade-offs between bit-depth and sample complexity [20, 41] for signal reconstruction, as they may vary for different problems when using specific algorithms. In contrast, we derive the Cramér-Rao bound for *parameter estimation* using quantized compressive random measurements, which provides an estimation-theoretic baseline for gauging the trade-off as well as benchmarking performances. Our CRB adds to existing literature of CRB calculations for 1-bit quantized single-tone frequency estimation [42] as well as for parameter estimation using compressive measurements [43].

1.4 Learning One-Hidden-Layer Neural Networks for Binary Classification

We then study the problem of learning neural-network-based models in the context of machine learning, i.e., we assume the training samples $(\mathbf{x}_i, y_i) \sim (\mathbf{x}, y)$, $i = 1, \dots, n$, are generated independently and identically distributed (i.i.d.) from a distribution based on a neural network model with the ground truth parameter $\mathbf{W}^* \in \mathbb{R}^{d \times K}$. The

output of an one-hidden-layer neural network can be written as

$$f(\mathbf{W}^*, \mathbf{x}_i) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{x}_i^\top \mathbf{w}_k^*), \quad (1.4)$$

where $\phi(\cdot)$ is the activation function, and \mathbf{w}_k^* denotes the k -th column of \mathbf{W}^* . The goal is to recover the underlying model parameter \mathbf{W}^* using the training samples. The non-linearity comes from two-fold. First the activation function considered is either Sigmoid $\phi(x) = \frac{1}{1+\exp(-x)}$ or ReLU $\phi(x) = \max(x, 0)$, both of which are nonlinear. Second with the binary classification setting we considered, y_i is quantized to discrete label $\{-1, 1\}$.

Neural networks have attracted a significant amount of research interest in recent years due to the success of deep neural networks [44] in practical domains such as computer vision and artificial intelligence [45–47]. Extensive studies have established the expressive power of neural networks. In particular, one-hidden-layer neural networks are sufficient to approximate any continuous function under certain conditions. However, the theoretical underpinnings behind this model remains mysterious to a large extent. This motivates us to study such a one-hidden-layer neural networks.

As with learning the linear model (1.1), an optimization problem (1.2) will be formulated for learning the neural-network-based model as well. Due to the non-linearity of the neural-network-based model, the optimization problem that needs to be solved is often non-convex. The intricate and unknown landscape of the non-convex objection function makes the problem much more difficult than the convex case.

Contributions

The main purpose of this study is to answer whether the true parameter can be recovered from its finite non-linear observations via solving the formulated non-convex optimization problem. As most machine learning tasks can be categorized as a classification problem or a regression problem, we considered both settings in this study, i.e.,

- *Regression*, where each sample $y \in \mathbb{R}$ is generated as

$$y = f(\mathbf{W}^*, \mathbf{x}).$$

This type of regression problem has been studied in various settings. In particular, [48] studied the single-neuron model under the Rectified Linear Unit (ReLU) activation, [7] studied the one-hidden-layer multi-neuron network model, and [49] studied a two-layer feedforward network with ReLU activations and identity mapping.

- *Classification*, where a label $y \in \{0, 1\}$ is drawn according to the conditional distribution

$$\mathbb{P}(y = 1 | \mathbf{x}) = f(\mathbf{W}^*, \mathbf{x}).$$

Such a problem has been studied in [50] when the network contains only a single neuron.

For both cases, previous studies attempted to recover \mathbf{W}^* , by minimizing an empirical loss function using the squared loss, i.e. $\min_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{W}, \mathbf{x}_i))^2$, given the training data. Two types of statistical guarantees were provided for such model recovery problems using the squared loss. More specifically, [7] showed that in the

local neighborhood of the ground truth \mathbf{W}^* , the *empirical* loss function is strongly convex for each *given* point under *independent* high probability event, which implies that *fresh samples* are required at *every* iteration for gradient descent to converge linearly with well-designed initializations. On the other hand, studies such as [50] established strong convexity in the entire local neighborhood of the ground truth in a uniform sense, so that resampling per iteration is not needed for gradient descent to have guaranteed linear convergence as long as it enters such a local neighborhood. Here, one weakness of the pointwise strong convexity in [7], compared to the uniform strong convexity in [50], is that independent fresh samples are required at each iteration to guarantee the linear convergence of gradient descent. Consequently, the sample complexity of [7] grows with respect to the recovery accuracy ϵ , typically with an extra factor of $\log(1/\epsilon)$ under linear convergence, which can be large when the desired accuracy is high. Therefore, the latter type of uniform strong convexity *without requiring per-iteration resampling* is much stronger and more desirable.

Considering the multi-neuron classification problem with either FCN or CNN. We show that the empirical risk function $L_n(\mathbf{W})$ is *uniformly* strongly convex in a local neighborhood of the ground truth \mathbf{W}^* , hence if initialized in this neighborhood, gradient descent converges linearly to a critical point (which we show to exist). Due to the nature of quantized labels here, the recovery of the ground truth is only up to certain statistical accuracy. Finally, we show that a tensor method provably provides an initialization in the neighborhood of the ground truth both for FCN and CNN.

The cross-entropy loss is much more challenging to analyze than the squared loss, e.g., its gradient and Hessian take much more complicated forms compared with the squared loss; moreover, it is hard to control the values of gradient and Hessian due

to the saturation phenomenon, i.e., when $f(\mathbf{W}, \mathbf{x})$ approaches 0 or 1. In order to establish the uniform local strong convexity property for the cross-entropy loss, we first show the *population* loss is smooth regarding to \mathbf{W}^* . Such a property was also established in [7] for the squared loss. However, considering the special form of Hessian under the cross-entropy loss, we need to apply Taylor’s approximation together with certain probabilistic upper bounds to control the value of Hessian, and obtain the smooth property. Network-specific quantities to capture the local geometry of the population loss at \mathbf{W}^* for FCN and CNN are derived, which imply that the geometry of CNN is more benign than FCN, corroborated by the numerical experiments. Beyond these two steps, the additional uniform concentration property of the Hessian is of key importance for us to obtain the uniform local strong convexity of the *empirical* loss. To show the uniform concentration of the Hessian, we successfully apply a type of covering argument. Different from the arguments in [50], which deal with the squared loss and are facilitated by certain nice assumptions on the activation functions, the cross-entropy loss is more difficult to apply the covering argument, e.g., both the gradient and Hessian no longer have a deterministic upper bound. Hence, we exploit the property of the sigmoid activation to show that the gradient and the Hessian of the cross-entropy loss are upper bounded with high probability in order to establish the uniform concentration property.

Related Work

Due to the scope, we focus on the most relevant literature on theoretical and algorithmic aspects of learning shallow neural networks via nonconvex optimization. The parameter recovery viewpoint is relevant to the success of nonconvex learning in

signal processing problems such as matrix completion, phase retrieval, blind deconvolution, dictionary learning and tensor decomposition [51–58], to name a few; see also the overview article [59]. The statistical model for data generation effectively removes worst-case instances and allows us to focus on average-case performance, which often possess much benign geometric properties that enable global convergence of simple local search algorithms.

The studies of one-hidden-layer network model can be further categorized into two classes, landscape analysis and model recovery. In the landscape analysis, it is known that if the network size is large enough compared to the data input, then there are no spurious local minima in the optimization landscape, and all local minima are global [60–63]. For the case with multiple neurons ($2 \leq K \leq d$) in the underparameterized setting, the work of Tian [64] studied the landscape of the population squared loss surface with ReLU activations. In particular, there exist spurious bad local minima in the optimization landscape [65,66] even at the population level. Zhong et. al. [7] provided several important geometric characterizations for the regression problem using a variety of activation functions and the squared loss.

In the model recovery problem, the number of neurons is smaller than the input dimension, and all the existing works discussed below assumed the squared loss and (sub-)Gaussian inputs. In the case with a single neuron ($K = 1$), [48] showed that gradient descent converges linearly when the activation function is ReLU, with a zero initialization, as long as the sample complexity is $O(d)$ for the regression problem. When the activation function is quadratic, [67] shows that randomly initialized gradient descent converges fast to the global optimum at a near-optimal sample complexity. On the other hand, [50] showed that when $\phi(\cdot)$ has bounded first, second and third

derivatives, there is no other critical points than the unique global minimum (within a constrained region of interest), and (projected) gradient descent converges linearly with an arbitrary initialization, as long as the sample complexity is $O(d \log^2 d)$ for the classification problem. Moreover, in the case with multiple neurons, [68] showed that projected gradient descent with a local initialization converges linearly for smooth activations with bounded second derivatives for the regression problem, [69] showed that gradient descent with tensor initialization converges linearly to a neighborhood of the ground truth using ReLU activations, and [70] showed the linear convergence of gradient descent with the spectral initialization using quadratic activations. For CNN with ReLU activations, [71] shows that gradient descent converges to the ground truth with random initialization for the population risk function based on the squared loss under Gaussian inputs. Moreover, [72] shows that gradient descent successfully learns a two-layer convolutional neural network despite the existence of bad local minima. From a technical perspective, our study differs from all the aforementioned work in that the cross entropy loss function we analyze has a very different form. Furthermore, we study the model recovery classification problem under the multi-neuron case, which has not been studied before.

1.5 Guaranteed Recovery of CNN with ReLU Activations

Moving beyond the classification setting, we further study a similar model recovery problem under the regression setting using the quadratic loss. In particular, we consider the training data that are generated by a one-hidden-layer convolutional neural network (CNN) with the ReLU activation function. Previous work [71, 73]

focused on analyzing the corresponding population risk function given by

$$L(\mathbf{w}) = \mathbb{E}_{\mathcal{D}} [\ell(\mathbf{w}; \mathbf{x}_i, y_i)], \quad (1.5)$$

where \mathcal{D} denotes the joint distribution of (\mathbf{x}, y) . They analyzed the performance of gradient descent that minimizes the population risk (1.5) with suitable assumptions on the joint distribution \mathcal{D} , and [71] further characterized the critical points of the population risk function. However, these studies considered only the first-order property about the population risk function. The second-order geometric property of the population risk function has not been explored, which is our interest here. We further wish to leverage the power of such second-order geometric properties for improving the previous results. Moreover, the performance of gradient descent on minimizing the empirical risk function has not been understood yet. Due to the non-smoothness of the ReLU activation, the landscape of the empirical risk function cannot be directly studied via its second-order property, i.e., the Hessian, and new methods need to be developed.

Contributions

We first provide a refined analysis of the landscape of the population risk function $L(\mathbf{w})$ under the Gaussian input assumption [71, 73], and show that it is strongly convex in a local neighborhood of the ground truth \mathbf{w}^* . Then, we show that gradient descent with random initialization converges to the global minimum \mathbf{w}^* at a linear rate. We further study the empirical risk function and show that such a function satisfies a regularity condition within a local neighborhood of \mathbf{w}^* . We then further establish that with a good initialization gradient descent converges linearly to the true weights \mathbf{w}^* .

Related Work

Besides the related work already mentioned in section 1.4, we here focus only on the studies of the model recovery problem with ReLU activation. [64] considered a single neuron neural network with ReLU activation, and showed that gradient descent with random initialization learns the true weights. Later, [49] showed that for a one-hidden-layer feed forward neural network with ReLU activation, SGD converges to the global minimum in two phases with small initialization weights. For the one-hidden-layer non-overlap convolutional neural network with ReLU activation, [71] studied the convergence of gradient descent with a fixed output layer and proved that gradient descent with random initialization can recover the weights exactly. Subsequently [73] generalized their result by not fixing the output layer, and more specifically they showed that when the output layer is to be learned, there exists spurious local minimum, but gradient descent still converges to the true weights with random initialization. These works have all been obtained under the assumption that input training data follow an i.i.d. Gaussian distribution. Moreover, all the convergence analysis of gradient descent methods is with regard to the population risk function.

In the finite sample regime, training the neural network is based on the empirical risk function. [48] studied the single-neuron network with ReLU activation in the high dimensional regime, and showed that projected gradient descent with initialization at the origin can recover the true weights up to some numerical constant under the assumption that the weights belong to a closed set and the input follows a Gaussian distribution.

The most related work to our study is [71], which studied only the population risk function. Our work captures the second-order landscape property of the population risk function, based on which we significantly improve the convergence rate of gradient descent from polynomial in [71] to linear. Inspired by the success of nonconvex learning in signal processing problems [51, 58, 74], we exploit the geometric property of the nonconvex empirical risk function and show the convergence of gradient descent with well-designed initialization.

1.6 Notations

Throughout this thesis, we use boldface letters to denote vectors and matrices, e.g. \mathbf{a} and \mathbf{A} . The Hermitian transpose of \mathbf{a} is denoted by \mathbf{a}^H , the transpose of \mathbf{a} is denoted by \mathbf{a}^\top , and $\|\mathbf{A}\|$, $\|\mathbf{A}\|_F$, $\text{Tr}(\mathbf{A})$ denote the spectral norm, the Frobenius norm, and the trace of the matrix \mathbf{A} , respectively. If \mathbf{A} is positive semidefinite (PSD), then $\mathbf{A} \succeq 0$. The identity matrix is denoted by \mathbf{I} . The gradient and the Hessian of a function $f(\mathbf{W})$ is denoted by $\nabla f(\mathbf{W})$ and $\nabla^2 f(\mathbf{W})$, respectively. We use c, C, C_1, \dots to denote constants whose values may vary from place to place. If necessary, we will introduce additional notations following the convention of notations in specific context in each chapter.

Chapter 2: Line Spectrum Estimation with Quantized Measurements

In this chapter, we consider the problem of recovering spectrally-sparse signals from heavy quantizations of their noisy complex-valued random linear measurements, e.g. only the quadrant information. We first derive the Cramér-Rao bound (CRB) for estimating multiple frequencies and their complex amplitudes assuming additive white Gaussian noise (AWGN) and the Lloyd's quantizer, using a fixed and deterministic compressed sensing (CS) measurement matrix. Furthermore, we use the derived CRB to answer the following question: given the same budget of bits, should we use more measurements (high sample complexity) with low bit-depth, or fewer measurements (low sample complexity) with high bit-depth? Moreover, we propose a novel atomic norm soft thresholding (AST) algorithm [29] to recover spectrally-sparse signals and estimate the frequencies from their 1-bit quantized measurements.

In this chapter we adopt the following notations. An indicator function for an event A is denoted as \mathbb{I}_A . Denote $\mathcal{T}(\mathbf{u}) \in \mathbb{C}^{n \times n}$ as the Hermitian Toeplitz matrix with \mathbf{u} as the first column. Define the inner product between two vectors \mathbf{a}, \mathbf{b} as $\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^H \mathbf{b}$. The cardinality of a set \mathcal{D} is defined as $|\mathcal{D}|$. If \mathbf{A} is positive semidefinite (PSD), then $\mathbf{A} \succeq 0$. $\Re(y)$ and $\Im(y)$ denote the real and imaginary part of a complex

number y , respectively. The expectation of a random variable a is written as $\mathbb{E}[a]$. Define \odot as entry-wise product.

2.1 Problem Formulation

Let $\mathbf{x}^* \in \mathbb{C}^n$ be a line spectrum signal, which is composed of a small number of spectral lines, defined as

$$\mathbf{x}^* = \sum_{k=1}^K c_k \mathbf{v}(f_k), \quad (2.1)$$

where K is the number of frequencies or level of sparsity, $c_k = A_k e^{j2\pi\phi_k} \in \mathbb{C}$ is the k th coefficient, $A_k > 0$ is the k th amplitude, $\phi_k \in [0, 1)$ is the k th normalized phase, $f_k \in [0, 1)$ is the k th frequency, and

$$\mathbf{v}(f) = [1 \quad e^{j2\pi f} \quad \dots \quad e^{j2\pi(n-1)f}]^T.$$

In CS, we acquire a set of random linear measurements of \mathbf{x}^* , contaminated by additive complex Gaussian noise, where each measurement is given as

$$z_i = \langle \mathbf{a}_i, \mathbf{x}^* \rangle + \sigma \epsilon_i, \quad i = 1, \dots, m, \quad (2.2)$$

where m is the number of measurements, $\mathbf{a}_i \in \mathbb{C}^n$'s are the measurement vectors composed of i.i.d. standard complex Gaussian entries $\mathcal{CN}(0, 1)$, σ is the noise level, and we further have i.i.d. $\epsilon_i \sim \mathcal{CN}(0, 1)$. In a vector notation, we write

$$\mathbf{z} = \mathbf{A} \mathbf{x}^* + \sigma \boldsymbol{\epsilon}, \quad (2.3)$$

where $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]^H \in \mathbb{C}^{m \times n}$ is the measurement matrix, $\boldsymbol{\epsilon} = [\epsilon_1, \epsilon_2, \dots, \epsilon_m]^T$, and $\mathbf{z} = [z_1, z_2, \dots, z_m]^T$. These measurements are then quantized into a finite number of bits for the ease of digital storage and processing. Denote $\mathcal{Q}(\cdot) : \mathbb{R} \mapsto \mathcal{D}$ as the quantizer that quantizes a real number into a finite alphabet \mathcal{D} , where the *bit depth*

is the smallest number of bits necessary to represent \mathcal{D} , i.e. $b^* = \min\{b \in \mathbb{Z}^+ : |\mathcal{D}| \leq 2^b\}$. The quantized measurements $\mathbf{y} = [y_1, y_2, \dots, y_m]^\top$ of \mathbf{z} are then denoted as

$$y_i = \mathcal{Q}(\Re(z_i)) + j\mathcal{Q}(\Im(z_i)), \quad i = 1, \dots, m, \quad (2.4)$$

where we apply the same quantizer \mathcal{Q} to both the real part and the imaginary part of the complex-valued measurement z_i . With slight abuse of notation, we denote the quantized measurements as

$$\mathbf{y} = \mathcal{Q}(\mathbf{z}). \quad (2.5)$$

Our goal is then to recover \mathbf{x}^* , and the set of frequencies $\mathbf{f} = \{f_k\}_{k=1}^K$, from the quantized measurements \mathbf{y} , possibly without a priori knowing the sparsity level K , and the form of the quantizer \mathcal{Q} .

Several choices of the quantizer are of special interest. At the extreme, we consider only knowing the quadrature information of z_i , where

$$Q(a) = \text{sign}(a), \quad a \in \mathbb{R}, \quad (2.6)$$

We refer to this quantizer as the *one-bit quantizer*, as only a single bit is used to quantize each real number.

More generally, we consider a quantizer $\mathcal{Q}(\cdot)$ that is fully characterized by the quantization intervals $\{[t_\ell, t_{\ell+1})\}_{\ell=1}^{|\mathcal{D}|-1}$, where $t_0 = -\infty$, $t_{|\mathcal{D}|} = \infty$, $\cup_{\ell=1}^{|\mathcal{D}|} [t_\ell, t_{\ell+1}) = \mathbb{R}$, as well as the representatives of each interval $\omega_\ell \in [t_\ell, t_{\ell+1})$, where

$$Q(a) = \omega_\ell, \quad \text{if } a \in [t_\ell, t_{\ell+1}). \quad (2.7)$$

For example, the Lloyd's quantizer [8] belongs to this form. The choice of the quantization scheme plays an important role in determining the performance of parameter estimation.

2.2 Cramer-Rao Bounds and Trade-offs

In this section, we study the effects of quantization on parameter estimation by deriving the Cramér-Rao bound assuming the quantizer, the sparsity level, and the noise level are known. In particular, the bounds are calculated for 1-bit and general quantizations, respectively, which are then used to study the trade-off between sample complexity and bit depths for a fixed bit budget.

To begin with, we assume the set of parameters, including the frequencies, amplitudes, and phases, given as $\boldsymbol{\kappa} = \{f_k, A_k, \phi_k\}_{k=1}^K \in \mathbb{R}^{3K}$, is deterministic but unknown, the measurement matrix \mathbf{A} is deterministic and known. Denote the probability mass function as $p(\mathbf{y}|\boldsymbol{\kappa})$, which is given as

$$p(\mathbf{y}|\boldsymbol{\kappa}) = \prod_{i=1}^m p(y_i|\boldsymbol{\kappa}) = \prod_{i=1}^m [p(\Re(y_i)|\boldsymbol{\kappa}) \cdot p(\Im(y_i)|\boldsymbol{\kappa})], \quad (2.8)$$

where the second equality follows from the fact that ϵ_i is proper. Moreover, let

$$p(\Re(y_i)|\boldsymbol{\kappa}) = \prod_{\omega \in \mathcal{D}} p_{\Re(y_i)}(\omega|\boldsymbol{\kappa})^{\mathbb{I}_{\{\Re(y_i)=\omega\}}} \quad (2.9)$$

$$p(\Im(y_i)|\boldsymbol{\kappa}) = \prod_{\omega \in \mathcal{D}} p_{\Im(y_i)}(\omega|\boldsymbol{\kappa})^{\mathbb{I}_{\{\Im(y_i)=\omega\}}} \quad (2.10)$$

be the probability mass function of $\Re(y)$ and $\Im(y)$, respectively. The Fisher Information Matrix (FIM), denoted by $\mathbf{I}(\boldsymbol{\kappa}) \in \mathbb{R}^{3K \times 3K}$, is given as

$$\mathbf{I}(\boldsymbol{\kappa}) = \mathbb{E} \left[\left(\frac{\partial \log p(\mathbf{y}|\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right) \left(\frac{\partial \log p(\mathbf{y}|\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top \right]. \quad (2.11)$$

Note that for any $1 \leq i, j \leq m$,

$$\begin{aligned} & \mathbb{E} \left[\left(\frac{\partial \log p(\Re(y_i)|\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right) \left(\frac{\partial \log p(\Im(y_j)|\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top \right] \\ &= \mathbb{E} \left[\frac{\partial \log p(\Re(y_i)|\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right] \cdot \mathbb{E} \left[\frac{\partial \log p(\Im(y_j)|\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right]^\top = \mathbf{0}, \end{aligned}$$

where the first equality follows from independence of $\Re(y_i)$ and $\Im(y_j)$, and the second equality follows from the fact

$$\begin{aligned} & \mathbb{E} \left[\left(\frac{\partial \log p(\Re(y_i) | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right) \right] \\ &= \mathbb{E} \left[\sum_{\omega \in \mathcal{D}} \frac{\mathbb{I}_{\{\Re(y_i) = \omega\}}}{p_{\Re(y_i)}(\omega | \boldsymbol{\kappa})} \frac{\partial p_{\Re(y_i)}(\omega | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right] \\ &= \sum_{\omega \in \mathcal{D}} \frac{\partial p_{\Re(y_i)}(\omega | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} = \frac{\partial (\sum_{\omega \in \mathcal{D}} p_{\Re(y_i)}(\omega | \boldsymbol{\kappa}))}{\partial \boldsymbol{\kappa}} = \mathbf{0}. \end{aligned}$$

Thus, plugging (2.8) into (2.11) all cross-terms will be zero and we have

$$\mathbf{I}(\boldsymbol{\kappa}) = \sum_{i=1}^m [\mathbf{I}_i^R(\boldsymbol{\kappa}) + \mathbf{I}_i^I(\boldsymbol{\kappa})], \quad (2.12)$$

where

$$\begin{aligned} \mathbf{I}_i^R(\boldsymbol{\kappa}) &= \mathbb{E} \left[\left(\frac{\partial \log p(\Re(y_i) | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right) \left(\frac{\partial \log p(\Re(y_i) | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top \right] \\ &= \sum_{\omega \in \mathcal{D}} \frac{1}{p_{\Re(y_i)}(\omega | \boldsymbol{\kappa})} \left(\frac{\partial p_{\Re(y_i)}(\omega | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right) \left(\frac{\partial p_{\Re(y_i)}(\omega | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top, \end{aligned}$$

and $\mathbf{I}_i^I(\boldsymbol{\kappa})$ can be given similarly by replacing $\Re(y_i)$ with $\Im(y_i)$.

The CRB for estimating $\boldsymbol{\kappa}$, is then given as $\text{CRB}(\boldsymbol{\kappa}) = \mathbf{I}(\boldsymbol{\kappa})^{-1}$, and the CRB for estimating the i th parameter in $\boldsymbol{\kappa}$, is given as $[\mathbf{I}(\boldsymbol{\kappa})^{-1}]_{i,i}$.

2.2.1 CRB for 1-Bit Quantization

Our goal is then to calculate the FIM in (2.12). We will explain in details the calculations for the 1-bit case. First, since $\Re(z_i) \sim \mathcal{N}(\Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle), \frac{1}{2}\sigma^2)$, then

$$p_{\Re(y_i)}(\omega | \boldsymbol{\kappa}) = \mathbb{P}(\omega \cdot \Re(z_i) > 0 | \boldsymbol{\kappa}) = \frac{1}{2} + \omega \cdot \Phi \left(\frac{\Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)}{\sigma} \right), \quad \omega = \pm 1, \quad (2.13)$$

where $\Phi(u) = \frac{1}{\sqrt{\pi}} \int_0^u e^{-t^2} dt$. Therefore, by the chain rule,

$$\begin{aligned} \frac{\partial p_{\Re(y_i)}(\omega | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} &= \frac{\omega}{\sigma} \Phi' \left(\frac{\Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)}{\sigma} \right) \frac{\partial \Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)}{\partial \boldsymbol{\kappa}} \\ &= \frac{\omega}{\sqrt{\pi}\sigma^2} \exp \left(-\frac{\Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)^2}{\sigma^2} \right) \frac{\partial \Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)}{\partial \boldsymbol{\kappa}}. \end{aligned} \quad (2.14)$$

As a short-hand notation, denote $s_i(\boldsymbol{\kappa}) = \Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)$ and $r_i(\boldsymbol{\kappa}) = \Im(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)$.

Plug (2.14) into $\mathbf{I}_i^R(\boldsymbol{\kappa})$, we have

$$\mathbf{I}_i^R(\boldsymbol{\kappa}) = \frac{4 \exp(-2s_i(\boldsymbol{\kappa})^2/\sigma^2)}{\pi\sigma^2 \left[1 - 4\Phi^2\left(\frac{s_i(\boldsymbol{\kappa})}{\sigma}\right)\right]} \begin{pmatrix} \frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \\ \frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \end{pmatrix}^\top,$$

and similarly,

$$\mathbf{I}_i^I(\boldsymbol{\kappa}) = \frac{4 \exp(-2r_i(\boldsymbol{\kappa})^2/\sigma^2)}{\pi\sigma^2 \left[1 - 4\Phi^2\left(\frac{r_i(\boldsymbol{\kappa})}{\sigma}\right)\right]} \begin{pmatrix} \frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \\ \frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \end{pmatrix}^\top.$$

As a remark, when $\sigma = 0$, the amplitude of the signal cannot be recovered from the 1-bit measurements due to scaling ambiguity, and the FIM becomes singular in this case. Therefore, our expressions for CRB is valid when $\sigma \neq 0$.

2.2.2 CRB for General Quantization

We now explain the calculation for a general quantization scheme. For $\omega_\ell \in \mathcal{D}$, and a corresponding interval $[t_\ell, t_{\ell+1})$, we have

$$\begin{aligned} p_{\Re(y_i)}(\omega_\ell | \boldsymbol{\kappa}) &= \mathbb{P}(\Re(z_i) \in [t_\ell, t_{\ell+1}) | \boldsymbol{\kappa}) \\ &= \int_{\frac{t_\ell - s_i(\boldsymbol{\kappa})}{\sigma}}^{\frac{t_{\ell+1} - s_i(\boldsymbol{\kappa})}{\sigma}} \frac{1}{\sqrt{\pi}} e^{-t^2} dt \\ &= \Phi\left(\frac{t_{\ell+1} - s_i(\boldsymbol{\kappa})}{\sigma}\right) - \Phi\left(\frac{t_\ell - s_i(\boldsymbol{\kappa})}{\sigma}\right), \end{aligned} \quad (2.15)$$

then, following similar arguments, we have

$$\frac{\partial p_{\Re(y_i)}(\omega_\ell | \boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} = \frac{1}{\sqrt{\pi\sigma^2}} \left[e^{-\frac{(t_{\ell+1} - s_i(\boldsymbol{\kappa}))^2}{\sigma^2}} - e^{-\frac{(t_\ell - s_i(\boldsymbol{\kappa}))^2}{\sigma^2}} \right] \frac{\partial \Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)}{\partial \boldsymbol{\kappa}}.$$

Therefore, define

$$\Gamma_i^R(\boldsymbol{\kappa}) = \sum_{\ell=1}^{|\mathcal{D}|-1} \frac{\left[e^{-\frac{(t_{\ell+1} - s_i(\boldsymbol{\kappa}))^2}{\sigma^2}} - e^{-\frac{(t_\ell - s_i(\boldsymbol{\kappa}))^2}{\sigma^2}} \right]^2}{\Phi\left(\frac{t_{\ell+1} - s_i(\boldsymbol{\kappa})}{\sigma}\right) - \Phi\left(\frac{t_\ell - s_i(\boldsymbol{\kappa})}{\sigma}\right)}, \quad (2.16)$$

and

$$\Gamma_i^I(\boldsymbol{\kappa}) = \sum_{\ell=1}^{|\mathcal{D}|-1} \frac{\left[e^{-\frac{(t_{\ell+1}-r_i(\boldsymbol{\kappa}))^2}{\sigma^2}} - e^{-\frac{(t_{\ell}-r_i(\boldsymbol{\kappa}))^2}{\sigma^2}} \right]^2}{\Phi\left(\frac{t_{\ell+1}-r_i(\boldsymbol{\kappa})}{\sigma}\right) - \Phi\left(\frac{t_{\ell}-r_i(\boldsymbol{\kappa})}{\sigma}\right)}, \quad (2.17)$$

we have the following theorem for the expression of FIM in light of our derivations in the previous subsection.

Theorem 1. *The Fisher information matrix $\mathbf{I}(\boldsymbol{\kappa})$ for estimating the unknown parameter $\boldsymbol{\kappa}$ is given as*

$$\mathbf{I}(\boldsymbol{\kappa}) = \frac{1}{\pi\sigma^2} \sum_{i=1}^m \left(\Gamma_i^R(\boldsymbol{\kappa}) \frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \left(\frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top + \Gamma_i^I(\boldsymbol{\kappa}) \frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \left(\frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top \right). \quad (2.18)$$

It is worth mentioning the FIM depends only on the quantization intervals, *not* the value of representatives. In contrast, the FIM using the unquantized measurements \mathbf{z} is given as

$$\mathbf{I}_{\text{unquantized}}(\boldsymbol{\kappa}) = \frac{2}{\sigma^2} \sum_{i=1}^m \frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \left(\frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top + \frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \left(\frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} \right)^\top.$$

It remains to evaluate $\frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}}$ and $\frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}}$. Following the Wirtinger calculus [75], we have $\frac{\partial s_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} = \frac{\partial \Re(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)}{\partial \boldsymbol{\kappa}} = \frac{1}{2} \Re(\mathbf{a}_i^H \frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\kappa}})$ and $\frac{\partial r_i(\boldsymbol{\kappa})}{\partial \boldsymbol{\kappa}} = \frac{\partial \Im(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)}{\partial \boldsymbol{\kappa}} = \frac{1}{2} \Im(\mathbf{a}_i^H \frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\kappa}})$. Define

$$\mathbf{g}(f) = \frac{\partial \mathbf{v}(f)}{\partial f} = [0, j2\pi e^{j2\pi f}, \dots, j2\pi(n-1)e^{j2\pi(n-1)f}]^\top.$$

Then, for each of the parameters in $\boldsymbol{\kappa}$, we have, for $k = 1, \dots, K$,

$$\frac{\partial \mathbf{x}^*}{\partial f_k} = c_k \mathbf{g}(f_k), \quad (2.19a)$$

$$\frac{\partial \mathbf{x}^*}{\partial A_k} = e^{j2\pi\phi_k} \mathbf{v}(f_k), \quad (2.19b)$$

$$\frac{\partial \mathbf{x}^*}{\partial \phi_k} = j2\pi c_k \mathbf{v}(f_k). \quad (2.19c)$$

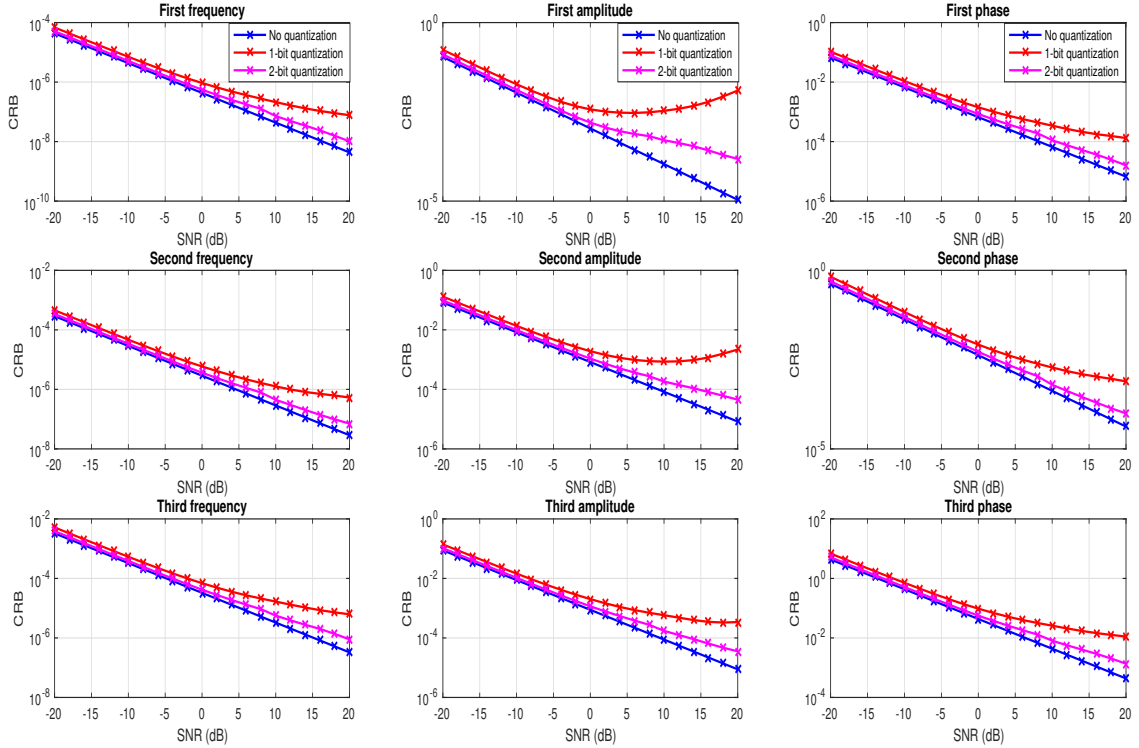


Figure 2.1: The CRB under different bit-depths with respect to SNR for a fixed number of measurements $m = 100$. Here, $n = 64$ and $K = 3$. Each row represents the CRB for estimating the frequency, amplitude and phase of one spectral atom.

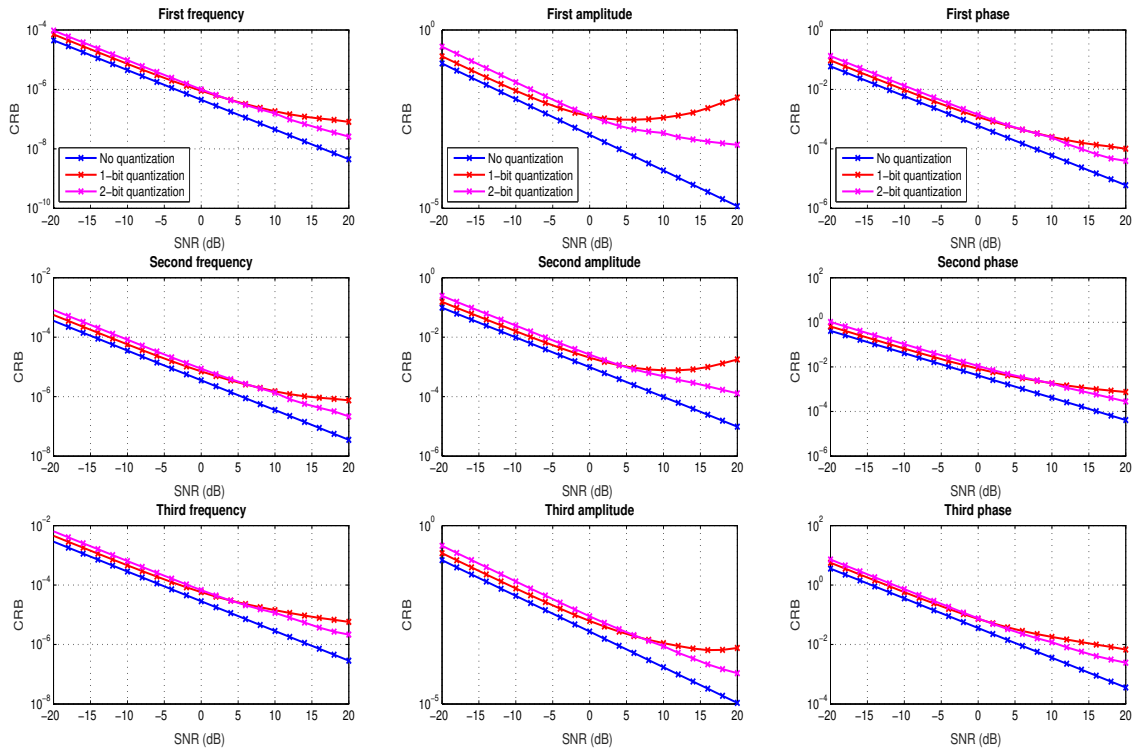


Figure 2.2: The CRB under different bit-depths with respect to SNR for a fixed number of bits $B = 100$. In this case, 2-bit quantization only has half the number of measurements of the 1-bit case. Here, $n = 64$ and $K = 3$. Each row represents the CRB for estimating the frequency, amplitude and phase of one spectral atom.

2.2.3 Numerical Evaluations of CRB

We now evaluate the CRB for 1-bit and 2-bit quantization schemes using the Lloyd’s quantizer, and compare it against the CRB without quantization. We generate a spectrally-sparse signal \mathbf{x}^* of length $n = 64$ with frequencies $f_1 = 0.3$, $f_2 = 0.325$, and $f_3 = 0.8$, and complex coefficients $c_1 = 0.4e^{j2\pi \cdot 0.1}$, $c_2 = 0.15e^{j2\pi \cdot 0.55}$, and $c_3 = 0.05e^{j2\pi \cdot 0.75}$, which are selected arbitrarily.

We first fix the number of measurements as $m = 100$, and generate a measurement matrix with complex standard i.i.d. Gaussian entries. Fig. 2.1 shows the CRB for estimating all parameters with respect to the SNR, where it is defined as $\text{SNR} = \|\mathbf{x}^*\|_2^2/\sigma^2$. It is evident that increasing the bit depth improves the performance. In the low SNR regime performance is *noise-limited*, and behaves similarly as if there was no quantization; in the high SNR regime, performance is *quantization-limited*, and experiences severe performance degeneration due to quantization.

In many situations, we cannot simultaneously have high sample complexity and high bit depth, but rather, our budget is constrained by the number of total bits, which is $B = m \cdot b^* = m \cdot \lceil \log_2 |\mathcal{D}| \rceil$. Therefore, it is useful to understand the trade-off between sample complexity and bit depth. Here, we use the CRB as a tool to compare the 1-bit and 2-bit quantization schemes. Fix the total number of bits as $B = 100$. In the 1-bit quantization scheme, we use a measurement matrix with $m = 100$ as generated earlier. In the 2-bit quantization scheme, we only use the first $m/2$ rows of the same measurement matrix. For comparison, we also plot the CRB assuming unquantized measurements using the same measurement matrix as the 1-bit case. Fig. 2.2 shows the CRB for estimating all parameters with respect to the SNR. It can be seen that in the low SNR regime, 1-bit quantization is preferred, as performance is

noise-limited, so higher sample complexity improves performance; in the high SNR regime, 2-bit quantization is preferred, as performance is quantization-limited, so higher bit depth improves performance. Our analysis is estimation-theoretic, and doesn't depend on the algorithm being adopted.

2.3 Atomic Norm Soft Thresholding for Quantized Spectral Compressed Sensing

It is well-known that maximum likelihood estimators approach the performance of CRB asymptotically at high SNR [21], however, their implementation requires exact knowledge of the likelihood function, which may not be available in certain applications. Therefore, in this section, we will develop estimators that do not require the knowledge of the quantization scheme using 1-bit measurements via atomic norm minimization [2]. We first provide the backgrounds on atomic norm for line spectrum estimation, and then describe the proposed algorithms for both the single vector case and the multiple vector case with performance guarantees.

2.3.1 Backgrounds on Atomic Norms

The atomic norm is originally proposed in [2] as a unified framework of convex regularizations for solving underdetermined linear inverse problems. Subsequently, [23–29] has tailored it to the estimation of spectrally-sparse signals. Most recently, [30] offered an extensive and comprehensive overview to atomic norm minimization as a canonical convex approach for super resolution.

For the single vector case, define the *atomic set* as

$$\mathcal{A}_s = \{e^{j\phi}\mathbf{v}(f) : f \in [0, 1), \phi \in [0, 2\pi)\},$$

then the atomic norm of a vector \mathbf{x} is given as

$$\|\mathbf{x}\|_{\mathcal{A}} := \inf\{t > 0 : \mathbf{x} \in t \cdot \text{conv}(\mathcal{A}_s)\} = \inf \left\{ \sum_i |\alpha_i| \mid \mathbf{x} = \sum_i \alpha_i \mathbf{v}(f_i) \right\}, \quad (2.20)$$

where $\text{conv}(\mathcal{A})$ denotes the convex hull of set \mathcal{A} . The atomic norm can be viewed as a continuous analog of the ℓ_1 norm over the continuous dictionary defined by the atomic set. Therefore, by promoting signals with small atomic norms, we encourage signals that can be expressed by a small number of spectral atoms. Appealingly, as shown in [23], it is possible to calculate $\|\mathbf{x}\|_{\mathcal{A}}$ using an equivalent semidefinite program, which can be computed efficiently using off-the-shelf solvers:

$$\|\mathbf{x}\|_{\mathcal{A}} = \min_{\mathbf{u} \in \mathbb{C}^n, w} \left\{ \frac{1}{2n} \text{Tr}(\mathcal{T}(\mathbf{u})) + \frac{w}{2} \mid \begin{bmatrix} \mathcal{T}(\mathbf{u}) & \mathbf{x} \\ \mathbf{x}^H & w \end{bmatrix} \succeq 0 \right\},$$

where $\mathcal{T}(\mathbf{u})$ denotes the Hermitian Toeplitz matrix with \mathbf{u} as the first column. The dual atomic norm $\|\cdot\|_{\mathcal{A}}^*$ for a vector $\mathbf{q} \in \mathbb{C}^n$, as will become useful later, is given as

$$\|\mathbf{q}\|_{\mathcal{A}}^* = \sup_{\|\mathbf{x}\|_{\mathcal{A}} \leq 1} \langle \mathbf{q}, \mathbf{x} \rangle_{\mathbb{R}} = \sup_{f \in [0,1]} |\mathbf{q}^H \mathbf{v}(f)|,$$

where the second equality follows from the fact the the extreme values are taken when \mathbf{x} is aligned with $\mathbf{v}(f)$ due to convexity. From the above equation it is clear that $\|\mathbf{q}\|_{\mathcal{A}}^*$ can be interpreted as the largest absolute value of a polynomial of $e^{j2\pi f}$, denoted as $Q(f) = |\mathbf{q}^H \mathbf{v}(f)|$.

2.3.2 Atomic Soft-Thresholding with Quantized Measurements

We first construct a surrogate signal from the quantized measurements as [37]

$$\mathbf{s} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i = \frac{1}{m} \mathbf{A}^H \mathbf{y} \in \mathbb{C}^n, \quad (2.21)$$

and use the following atomic norm soft-thresholding (AST) algorithm to estimate the signal \mathbf{x} ,

$$\hat{\mathbf{x}} = \underset{\mathbf{x} \in \mathbb{C}^n}{\text{argmin}} \frac{1}{2} \|\mathbf{x} - \mathbf{s}\|_2^2 + \tau \|\mathbf{x}\|_{\mathcal{A}}, \quad (2.22)$$

which is the proximal mapping of the surrogate signal \mathbf{s} with respect to the atomic norm, where $\tau > 0$ is a regularization parameter. One appealing feature of atomic norm minimization is that the set of frequencies can be recovered via the dual polynomial approach [29]. Namely, denote the dual variable as $\hat{\mathbf{q}} = (\mathbf{s} - \hat{\mathbf{x}})/\tau$, and $Q(f) = |\hat{\mathbf{q}}^H \mathbf{v}(f)|$. Then the set of frequencies can be localized as $\hat{\mathcal{F}} = \{f : Q(f) = 1\}$. We refer interested readers to the details in [23]. Alternatively, the frequencies can be localized via performing conventional subspace methods using the estimated signal.

2.3.3 Performance Guarantees

In this section, we develop performance guarantees of the proposed AST algorithm under 1-bit quantization in the single vector case using the sign quantizer in (2.6). Note that in this case, it can be seen that \mathbf{s} in (2.21) is an unbiased estimator of \mathbf{x}^* up to a scaling difference, i.e.

$$\mathbb{E}[\mathbf{s}] = \lambda \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2},$$

where

$$\lambda = \frac{2\|\mathbf{x}^*\|_2}{\sqrt{\pi(\sigma^2 + \|\mathbf{x}^*\|_2^2)}} = \frac{2}{\sqrt{\pi(1/\text{SNR} + 1)}} \quad (2.23)$$

depends on the SNR before quantization $\text{SNR} = \|\mathbf{x}^*\|_2^2/\sigma^2$. To illustrate, Fig. 2.3 depicts λ as a function of SNR, which is a monotonically increasing function with respect to SNR and approaches to the limit $2/\sqrt{\pi}$ as SNR goes to infinity.

Without loss of generality, we assume $\|\mathbf{x}^*\|_2 = 1$. The performance of AST relies critically on the separation condition, which is defined as the minimum distance between distinct frequencies,

$$\Delta = \min_{k \neq j} |f_k - f_j| \geq \frac{4}{n}, \quad (2.24)$$

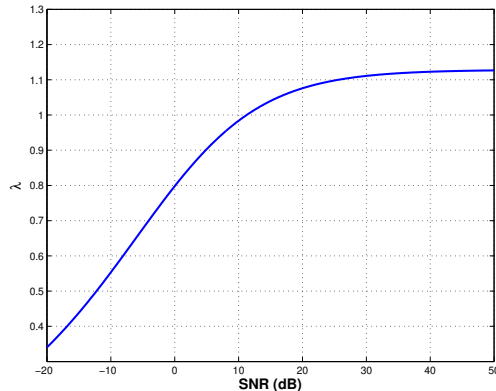


Figure 2.3: The value of λ with respect to SNR before quantization.

where $|f_k - f_j|$ is evaluated as the wrap-around difference on the unit modulus. Under the separation condition, we have the performance guarantee of the proposed algorithm in (2.22), stated below.

Theorem 2. *Set $\tau := \eta \sqrt{n \log n / m}$ for some constant $\eta \geq 1$. Under the separation condition, the solution $\hat{\mathbf{x}}$ satisfies*

$$\left\| \frac{\hat{\mathbf{x}}}{\lambda} - \mathbf{x}^* \right\|_2 \lesssim \frac{1}{\lambda} \sqrt{\frac{K \log n}{m}}$$

with high probability.

The proof of Theorem 2 can be found in Appendix A.1. Theorem 2 suggests that the proposed algorithm accurately recovers the signal as soon as m is on the order of $K \log n$, which is order-wise near-optimal, since at least an order of $K \log(n/K)$ measurements are needed in order to recover a sparse signal in the DFT basis [35]. Moreover, the theorem also suggests that the normalized reconstruction error is inverse proportional to λ , which plays the role of SNR *after quantization* and is a

nonlinear function of the SNR before quantization. In the low SNR regime, λ scales as $1/\sqrt{\text{SNR}}$, and the performance is comparable to that using unquantized measurements. However, in the high SNR regime, there is a saturation phenomenon, as evidenced by Fig. 2.3, and the performance does not improve as much with we increase SNR, which is also corroborated by numerical simulation in Section 2.5. These results are qualitatively in line with existing work on one-bit CS [35].

Remark: More generally, Theorem 2 can be extended to the generalized linear model following similar strategies in [34], as long as the 1-bit measurements y_i 's are i.i.d. and satisfy $\mathbb{E}[y_i|\mathbf{a}_i] = g(\langle \mathbf{a}_i, \mathbf{x}^* \rangle)$ for some link function $g(\cdot)$, and accordingly $\lambda = \mathbb{E}[g(\theta)\theta^H]$ where the expectation is taken with respect to $\theta \sim \mathcal{CN}(0, 1)$. This allows us to model other complex quantization schemes with non-Gaussian noise.

2.4 Extension to the Multiple Vector Case

In many applications, we encounter an ensemble of line spectrum signals, where each signal $\mathbf{x}_t \in \mathbb{C}^n$ contains a linear combination of spectral lines with the same set of frequencies \mathcal{F} , but with varying amplitudes, given as

$$\mathbf{x}_t^* = \sum_{k=1}^K c_{k,t} \mathbf{v}(f_k), \quad 1 \leq t \leq T,$$

where $c_{k,t} \in \mathbb{C}$, and T is the number of snapshots. Denote $\mathbf{X}^* = [\mathbf{x}_1^*, \mathbf{x}_2^*, \dots, \mathbf{x}_T^*] \in \mathbb{C}^{n \times T}$ as the signal ensemble. Similar to (2.3), the CS measurement of each snapshot is given as

$$\mathbf{z}_t = \mathbf{A} \mathbf{x}_t^* + \sigma \boldsymbol{\epsilon}_t, \quad (2.25)$$

where $\boldsymbol{\epsilon}_t = [\epsilon_{1,t}, \epsilon_{2,t}, \dots, \epsilon_{m,t}]^T$ contains i.i.d. standard complex Gaussian $\mathcal{CN}(0, 1)$ entries. Similar to (2.5), the quantized measurements of each \mathbf{z}_t is then given as

$$\mathbf{y}_t = \mathcal{Q}(\mathbf{z}_t). \quad (2.26)$$

Denote $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T]$ and $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T]$ as the unquantized measurement ensemble and the quantized measurement ensemble, respectively. Our goal is then to recover \mathbf{X}^* and the set of frequencies from \mathbf{Y} , without assuming the knowledge of the sparsity level and the quantizer. The presence of multiple vectors can significantly improve the accuracy of frequency estimation.

It is possible to extend the atomic norm formulation to the multiple vector case [25]. Define the atomic set as

$$\mathcal{A}_m = \{ \mathbf{A}(f, \mathbf{b}) = \mathbf{v}(f) \mathbf{b} \mid f \in (0, 1], \mathbf{b} \in \mathbb{C}^{1 \times T}, \|\mathbf{b}\|_2 = 1 \},$$

then the atomic norm is defined as

$$\begin{aligned} \|\mathbf{X}\|_{\mathcal{A}} &= \inf \{ t > 0 : \mathbf{X} \in t \cdot \text{conv}(\mathcal{A}_m) \} \\ &= \inf \left\{ \sum_k |c_k| \mid \mathbf{X} = \sum_k c_k \mathbf{A}(f_k, \mathbf{b}_k) \right\}, \end{aligned}$$

which can be computed similarly via solving the following semidefinite program [25]:

$$\|\mathbf{X}\|_{\mathcal{A}} = \min_{\mathbf{u} \in \mathbb{C}^n, \mathbf{W} \in \mathbb{C}^{T \times T}} \left\{ \frac{1}{2n} \text{Tr}(\mathcal{T}(\mathbf{u})) + \frac{1}{2} \text{Tr}(\mathbf{W}) \mid \begin{bmatrix} \mathcal{T}(\mathbf{u}) & \mathbf{X} \\ \mathbf{X}^H & \mathbf{W} \end{bmatrix} \succeq 0 \right\}.$$

The dual norm for some $\mathbf{Q} \in \mathbb{C}^{n \times T}$ is given as

$$\|\mathbf{Q}\|_{\mathcal{A}}^* = \sup_{\|\mathbf{X}\|_{\mathcal{A}} \leq 1} \langle \mathbf{Q}, \mathbf{X} \rangle_{\mathbb{R}} = \sup_{f \in [0, 1]} \|\mathbf{Q}^H \mathbf{v}(f)\|_2,$$

which is the largest absolute value of the polynomial $Q(f) = \|\mathbf{Q}^H \mathbf{v}(f)\|_2$.

For reconstruction, we construct the surrogate signal ensemble from the quantized measurement ensemble \mathbf{Y} as

$$\mathbf{S} = \frac{1}{m} \mathbf{A}^H \mathbf{Y} \in \mathbb{C}^{n \times T}, \quad (2.27)$$

and use the following atomic norm soft-thresholding (AST) algorithm to estimate the signal ensemble \mathbf{X} ,

$$\hat{\mathbf{X}} = \operatorname{argmin}_{\mathbf{X} \in \mathbb{C}^{n \times T}} \|\mathbf{X} - \mathbf{S}\|_F^2 + \tau_T \|\mathbf{X}\|_{\mathcal{A}}, \quad (2.28)$$

where $\tau_T > 0$ is a regularization parameter. Moreover, define $\hat{\mathbf{Q}} = (\mathbf{S} - \hat{\mathbf{X}})/\tau_T$, and $Q(f) = \|\hat{\mathbf{Q}}^H \mathbf{v}(f)\|_2$. Then the set of frequencies can be localized as $\hat{\mathcal{F}} = \{f : Q(f) = 1\}$. Alternatively, the frequencies can be localized via performing conventional subspace methods using the estimated snapshots.

2.5 Numerical Experiments

In this section, we conduct numerical experiments to evaluate the performance of the proposed AST algorithms for parameter estimation using quantized compressive measurements in both the single vector case and the multiple vector case. For implementation of the AST algorithms, we used the CVX toolbox [76]. There're several other fast solvers developed for atomic norm minimization that are more scalable to large problems, including ADMM [25, 29], ADCG [77], and CoGent [78], to name a few.

2.5.1 Single Vector Case

Let $n = 64$ and $K = 3$. The set of frequencies is located at $\mathbf{f} = \{0.3, 0.325, 0.8\}$, where the first two frequencies are separated barely more than $1/n$, the Rayleigh

limit. The number of bits is set as $m = 1000$, where the measurement vectors are generated with i.i.d. $\mathcal{CN}(0, 1)$ entries. The measurements are quantized according to (2.6). Fig. 2.4 shows the amplitude of the constructed dual polynomial by solving (2.22), where its peaks can be used to localize the frequencies. It can be seen that it matches accurately with the ground truth.

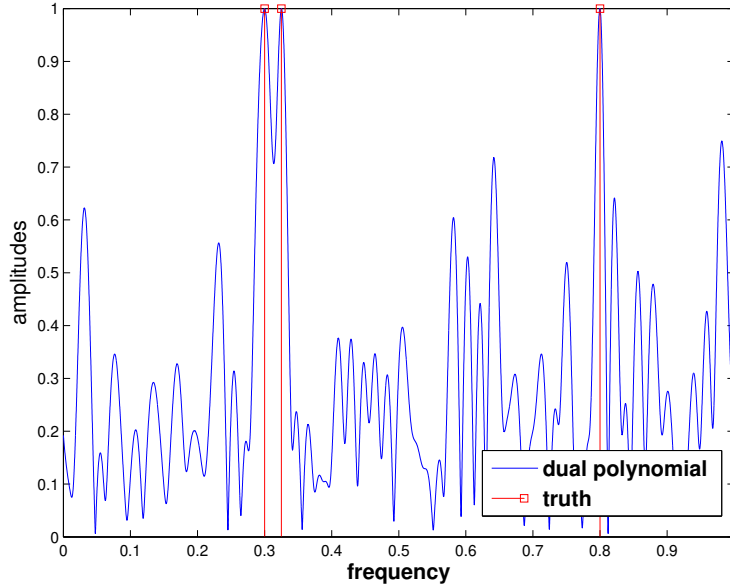


Figure 2.4: Frequency localization via peaks of the dual polynomial, superimposed on the ground truth.

Next, we compare the performance of signal reconstruction using atomic norm with unquantized measurements \mathbf{z} , by running the algorithm:

$$\hat{\mathbf{x}}_{\text{UQ}} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{C}^n} \frac{1}{2} \|\mathbf{z} - \mathbf{A}\mathbf{x}\|_2^2 + \tilde{\tau} \|\mathbf{x}\|_{\mathcal{A}},$$

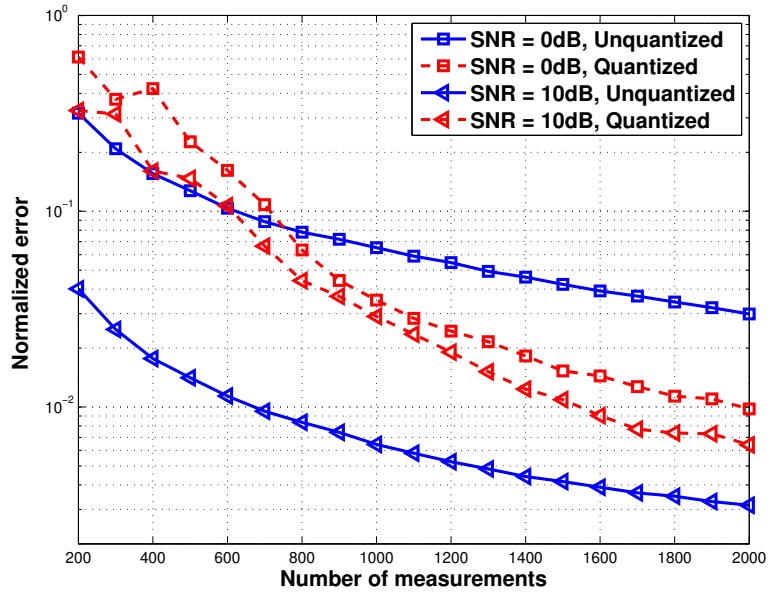


Figure 2.5: Normalized reconstruction error with respect to the number of measurements at different SNRs with or without quantization.

where $\tilde{\tau}$ is a properly tuned regularization parameter. The normalized reconstruction error is defined as $\sin^2(\angle \hat{\mathbf{x}}, \mathbf{x}) = 1 - |\langle \hat{\mathbf{x}}, \mathbf{x}^* \rangle|^2 / (\|\hat{\mathbf{x}}\|_2^2 \|\mathbf{x}^*\|_2^2)$, where $\hat{\mathbf{x}}$ is the reconstructed signal using either algorithm. Fig. 2.5 shows the normalized reconstruction error at different SNRs with comparisons to that using the quantized measurements and the AST algorithm (2.22), where SNR is defined again as $\text{SNR} = \|\mathbf{x}^*\|_2^2 / \sigma^2$. It can be seen that the reconstruction accuracy improves as we increase the SNR as well as the number of measurements, validating the theoretical analysis. In particular, at low SNR, using quantized measurements can potentially achieve better reconstruction quality with much fewer measurement budgets in bits. It can also be seen that

improving the SNR before quantization does not have as strong impact as for the unquantized case.

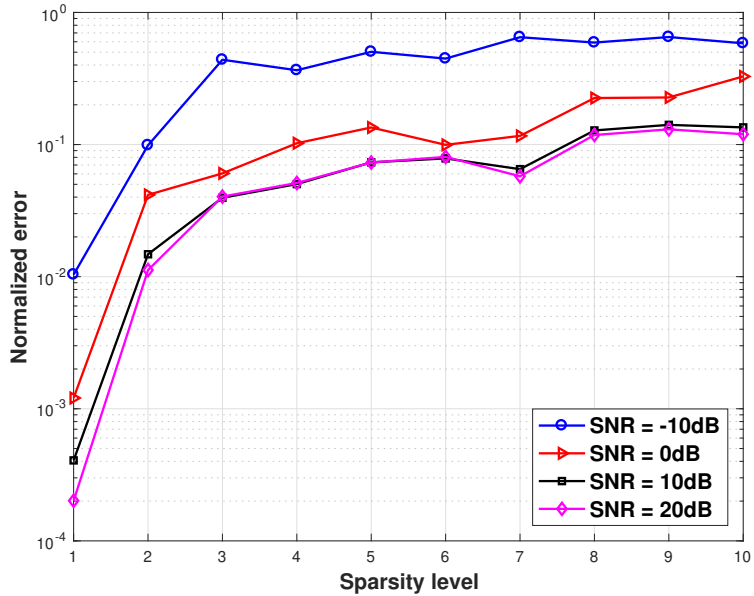


Figure 2.6: Normalized reconstruction error with respect to the spectral sparsity level at different SNRs before quantization.

Next, we examine the performance of the proposed algorithm as a function of the spectral sparsity level. Fix $n = 64$ and $m = 1000$. At each run, we randomly generate K different frequencies that satisfy the separate condition. Fig. 2.6 shows the normalized reconstruction error as a function of the sparsity level at various SNR, averaged over 200 Monte Carlo simulations. It can be seen that the reconstruction error is higher when the spectral sparsity level is higher, and the SNR is lower. Moreover, it can be seen that the reconstruction error stops to decrease when the

SNR is relatively high, indicating a saturation effect due to quantization, as predicted by our theory.

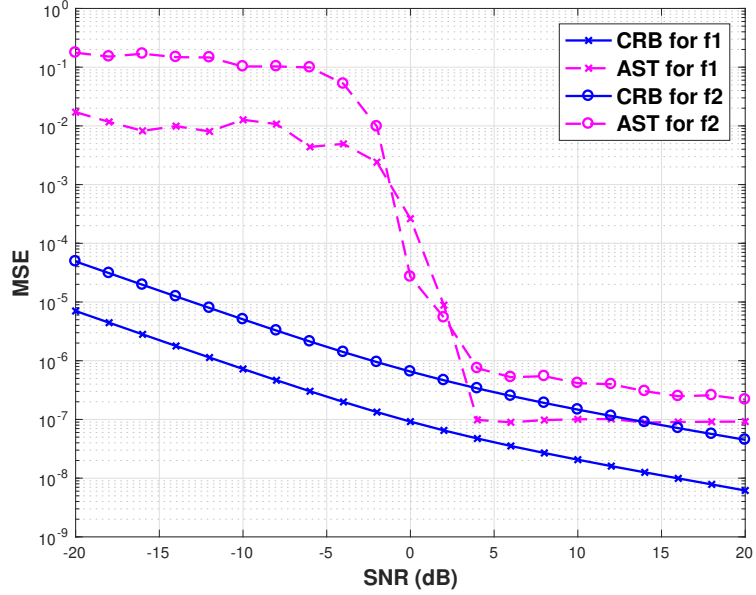


Figure 2.7: Mean square error of frequency localization with respect to SNR using 1-bit measurements, CRB is provided as a benchmark: (a) first frequency; (b) second frequency.

We further compare the performance of frequency localization using the proposed algorithm with the CRB. Fix $n = 64$ and $m = 1000$. We generate the ground signal with frequencies $f_1 = 0.3$, $f_2 = 0.325$ and amplitudes $c_1 = 0.4e^{j2\pi \cdot 0.1}$, $c_2 = 0.15e^{j2\pi \cdot 0.55}$. Fig. 2.7 shows the average mean squared error for each frequency over 200 Monte Carlo simulations, against the corresponding CRB calculated using the formulas in Section 2.2. The frequencies are estimated by using the MATLAB function `rootmusic` by assuming the correct model order, that is $K = 2$. The performance of the proposed

algorithm exhibits a threshold effect where it approaches that of CRB as soon as SNR is large enough. However, further increasing the SNR doesn't seem to improve the performance, which coincides with the saturation effect discussed earlier.

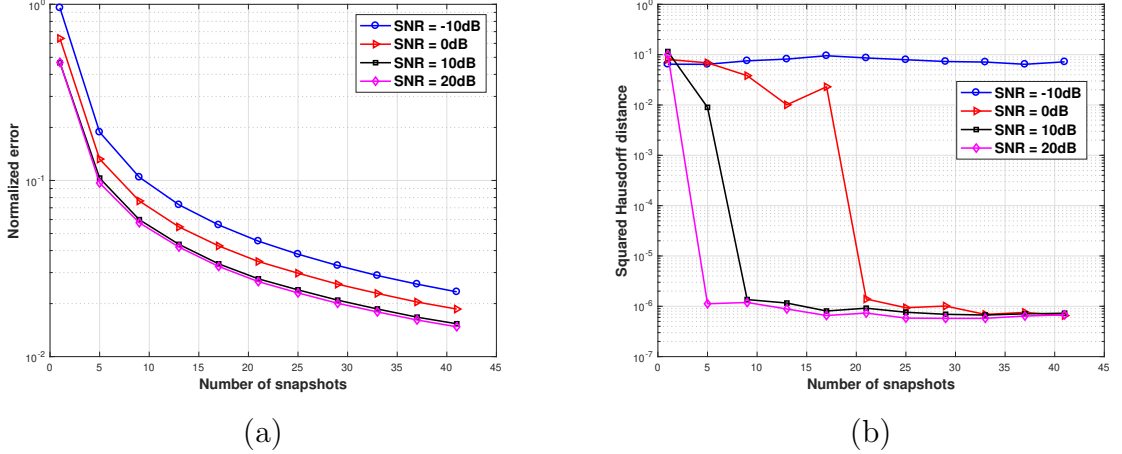


Figure 2.8: Performance with respect to the number of snapshots at different SNRs using 1-bit measurements: (a) signal reconstruction error; (b) frequency estimation error measured in Hausdorff distance.

2.5.2 Multiple Vector Case

We evaluate the performance of the AST algorithm (2.28) in the multiple vector case. We follow the same setup as Fig. 2.4, where $n = 64$, the set of frequencies $\mathbf{f} = \{0.3, 0.325, 0.8\}$, and the number of measurements for each snapshot is $m = 50$. The coefficients of each snapshot in \mathbf{X} is generated independently using the standard complex Gaussian distribution. The SNR *per snapshot* is defined as $\text{SNR} = \|\mathbf{X}\|_{\mathbb{F}}^2 / (T\sigma^2)$, where T is the number of snapshots. We set the regularization

parameter $\tau_T = \sqrt{n \log n / (10 \cdot mT)}$ in the experiment. The normalized reconstruction error is defined as $\sin^2(\angle \mathbf{X}; \hat{\mathbf{X}})$, where $\hat{\mathbf{X}}$ is the recovered signal containing multiple snapshots, and \angle denotes the angle between the subspace spanned by \mathbf{X} and $\hat{\mathbf{X}}$. Once $\hat{\mathbf{X}}$ is obtained, we estimate the frequencies by using the MATLAB function `rootmusic` by assuming the correct model order, that is $K = 3$. The accuracy of frequency estimation is evaluated by examining the Hausdorff distance between the recovered frequencies $\hat{\mathbf{f}}$ and the ground truth \mathbf{f} as

$$d_H(\mathbf{f}, \hat{\mathbf{f}}) = \max \left\{ \sup_{f \in \mathbf{f}} \inf_{\hat{f} \in \hat{\mathbf{f}}} \|f - \hat{f}\|_2, \sup_{\hat{f} \in \hat{\mathbf{f}}} \inf_{f \in \mathbf{f}} \|f - \hat{f}\|_2 \right\}.$$

Fig. 2.8 shows the recovery performance with respect to the number of snapshots at different SNRs, averaged over 50 Monte Carlo simulations, where (a) depicts the normalized reconstruction error, and (b) depicts the squared Hausdorff distance. At a fixed SNR, it can be seen that both the normalized reconstruction error and frequency estimation error reduce, highlighting the benefit of having multiple snapshots. In particular, having multiple snapshots allows better frequency recovery once the number of snapshots is large enough. Moreover, performance improves as we increase the SNR.

Chapter 3: Learning One-Hidden-Layer Neural Network for Binary Classification

In this chapter, we study model recovery for data classification, where the training labels are generated from a one-hidden-layer neural network with sigmoid activations, and the goal is to recover the weights of the neural network. Considering the multi-neuron classification problem with either FCN or CNN, we first analyze the landscape of the cross-entropy loss and show that gradient descent converges linearly to a critical point, which is shown to exist. Due to the quantized nature of labels, the recovery of the ground truth is only up to certain statistical accuracy. At last, we show that a tensor method provably provides an initialization in the neighborhood of the ground truth both for FCN and CNN.

We adopt additional notations in this chapter as follows. Denote $\|\cdot\|_{\psi_1}$ as the sub-exponential norm of a random variable. For nonnegative functions $f(x)$ and $g(x)$, $f(x) = O(g(x))$ means there exist positive constants c and a such that $f(x) \leq cg(x)$ for all $x \geq a$; $f(x) = \Omega(g(x))$ means there exist positive constants c and a such that $f(x) \geq cg(x)$ for all $x \geq a$.

3.1 Problem Formulation

We consider two popular types of one-hidden-layer nonlinear neural networks illustrated in Fig. 3.1, i.e., a FCN [7] and a non-overlapping CNN [71]. For both cases, we let $\mathbf{x} \in \mathbb{R}^d$ be the input, $K \geq 1$ be the number of neurons, and the activation function be the sigmoid function

$$\phi(x) = \frac{1}{1 + \exp(-x)}.$$

- *FCN*: the network parameter is $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K] \in \mathbb{R}^{d \times K}$, and

$$H_{\text{FCN}}(\mathbf{W}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}_k^\top \mathbf{x}). \quad (3.1)$$

- *Non-overlapping CNN*: for simplicity we let $d = m \cdot K$ for some integers m .

Let $\mathbf{w} \in \mathbb{R}^m$ be the network parameter, and the k th stride of \mathbf{x} be given as

$$\mathbf{x}^{(k)} = [x_{m(k-1)+1}, \dots, x_{m \cdot k}]^\top \in \mathbb{R}^m. \text{ Then,}$$

$$H_{\text{CNN}}(\mathbf{w}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}^\top \mathbf{x}^{(k)}). \quad (3.2)$$

The non-overlapping CNN model can be viewed as a highly structured instance of the FCN, where the weight matrix can be written as:

$$\mathbf{W}_{\text{CNN}} = \begin{bmatrix} \mathbf{w} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{w} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{w} \end{bmatrix} \in \mathbb{R}^{d \times K}.$$

In a model recovery setting, we are given n training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim (\mathbf{x}, y)$ that are drawn i.i.d. from certain distribution regarding the ground truth network parameter \mathbf{W}^* (or resp. \mathbf{w}^* for CNN). Suppose the network input $\mathbf{x} \in \mathbb{R}^d$ is drawn from a standard Gaussian distribution $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. This assumption has been used

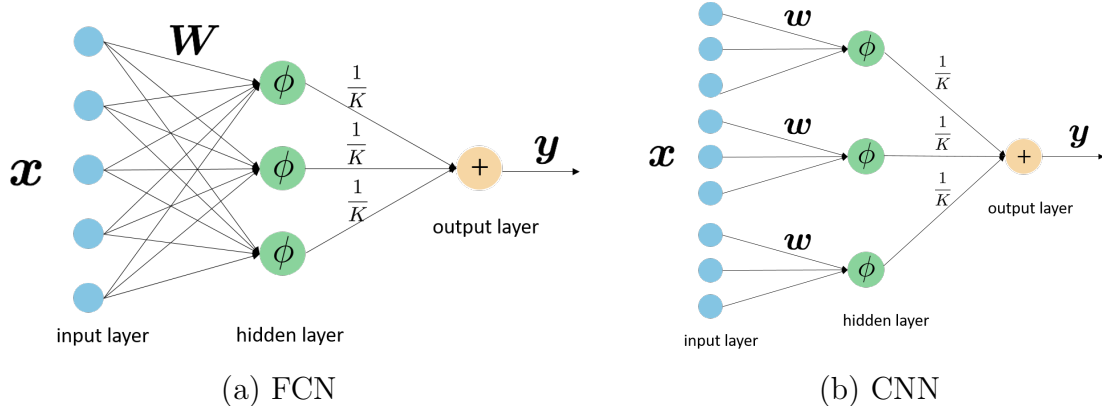


Figure 3.1: Illustration of two types of one-hidden-layer neural networks considered in this Chapter: (a) a fully-connected network (FCN); (b) a non-overlapping convolutional neural network (CNN).

a lot in previous literature [48,68,71,72], to name a few. Then, conditioned on $\mathbf{x} \in \mathbb{R}^d$, the output y is mapped to $\{0, 1\}$ via the output of the neural network, i.e.,

$$\mathbb{P}(y = 1|\mathbf{x}) = H(\mathbf{W}^*, \mathbf{x}). \quad (3.3)$$

Our goal is to recover the network parameter, i.e., \mathbf{W}^* , via minimizing the following empirical loss function:

$$f_n(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{W}; \mathbf{x}_i), \quad (3.4)$$

where $\ell(\mathbf{W}; \mathbf{x}) := \ell(\mathbf{W}; \mathbf{x}, y)$ is the cross-entropy loss function, i.e.,

$$\ell(\mathbf{W}; \mathbf{x}) = -y \cdot \log(H(\mathbf{W}, \mathbf{x})) - (1 - y) \cdot \log(1 - H(\mathbf{W}, \mathbf{x})), \quad (3.5)$$

where $H(\mathbf{W}, \mathbf{x})$ can subsume either H_{FCN} or H_{CNN} .

3.2 Gradient Descent and its Performance Guarantee

To estimate the network parameter \mathbf{W}^* , since (3.4) is a highly nonconvex function, vanilla gradient descent with an arbitrary initialization may get stuck at local minima. Therefore, we implement gradient descent (GD) with a well-designed initialization scheme that is described in details in Section 3.3. In this section, we focus on the performance of the local update rule

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla f_n(\mathbf{W}_t),$$

where η is the constant step size. The algorithm is summarized in Algorithm 1.

Algorithm 1 Gradient Descent (GD)

Input: Training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, step size η , iteration T

Initialization: $\mathbf{W}_0 \leftarrow \text{INITIALIZATION}(\{(\mathbf{x}_i, y_i)\}_{i=1}^n)$

Gradient Descent: for $t = 0, 1, \dots, T - 1$

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla f_n(\mathbf{W}_t).$$

Output: \mathbf{W}_T

Note that throughout the execution of GD, the same set of training samples is used which is the standard implementation of gradient descent. Consequently the analysis is challenging due to the statistical dependence of the iterates with the data.

3.2.1 Uniform local strong convexity

We first characterize the local strong convexity of $f_n(\cdot)$ in a neighborhood of the ground truth. We use the Euclidean ball to denote the local neighborhood of \mathbf{W}^* for

FCN or of \mathbf{w}^* for CNN.

$$\mathbb{B}(\mathbf{W}^*, r) = \{\mathbf{W} \in \mathbb{R}^{d \times K} : \|\mathbf{W} - \mathbf{W}^*\|_F \leq r\}, \quad (3.6a)$$

$$\mathbb{B}(\mathbf{w}^*, r) = \{\mathbf{w} \in \mathbb{R}^m : \|\mathbf{w} - \mathbf{w}^*\|_2 \leq r\}, \quad (3.6b)$$

where r is the radius of the ball. With slight abuse of notations, we will drop the subscript FCN or CNN for simplicity, whenever it is clear from the context that the result is for FCN when the argument is $\mathbf{W} \in \mathbb{R}^{d \times K}$ and for CNN when the argument is $\mathbf{w} \in \mathbb{R}^m$. Further, $\sigma_i(\mathbf{W})$ denotes the i -th largest singular value of \mathbf{W}^* . Let the condition number be $\kappa = \sigma_1/\sigma_K$, and $\lambda = \prod_{i=1}^K (\sigma_i/\sigma_K)$. Moreover, we introduce an important quantity $\rho(\sigma)$ regarding $\phi(z)$, the sigmoid activation function, that captures the geometric properties of the loss function for neural networks (3.1) and (3.2).

Definition 1 (Key quantity for FCN). *Let $z \sim \mathcal{N}(0, 1)$ and define $\alpha_q(\sigma) = \mathbb{E}[\phi'(\sigma \cdot z)z^q]$, $\forall q \in \{0, 1, 2\}$, and $\beta_q(\sigma) = \mathbb{E}[\phi'(\sigma \cdot z)^2 z^q]$, $\forall q \in \{0, 2\}$. Define $\rho_{\text{FCN}}(\sigma)$ as*

$$\rho_{\text{FCN}}(\sigma) = \min \{\beta_0(\sigma) - \alpha_0^2(\sigma), \beta_2(\sigma) - \alpha_2^2(\sigma)\} - \alpha_1^2(\sigma).$$

Definition 2 (Key quantity for CNN). *Let $z \sim \mathcal{N}(0, \sigma^2)$ and define $\rho_{\text{CNN}}(\sigma)$ as*

$$\rho_{\text{CNN}}(\sigma) = \min \left\{ \mathbb{E}[(\phi'(z)z)^2], \mathbb{E}[\phi'(z)^2] \right\}.$$

Note that Definition 1 for FCN is different from that in [7, Property 3.2] but consistent with [7, Lemma D.4] which removes the third term in [7, Property 3.2]. For the activation function considered in this Chapter, the first two terms suffice. Definition 2 for CNN is a newly distilled quantity in this Chapter tailored to the special structure of CNN.

The quantity $\rho(\sigma)$ plays an important role in the following theorem which guarantees the Hessian of the empirical risk function in the local neighborhood of the ground truth is positive definite with high probability for both FCN and CNN.

Theorem 3 (Local Strong Convexity). *Consider the classification model with FCN (3.1) or CNN (3.2) and the sigmoid activation function.*

- For FCN, assume $\|\mathbf{w}_k^*\|_2 \leq 1$ for all k . There exist constants c_1 and c_2 such that as soon as sample size

$$n_{\text{FCN}} \geq c_1 \cdot dK^5 \log^2 d \cdot \left(\frac{\kappa^2 \lambda}{\rho_{\text{FCN}}(\sigma_K)} \right)^2,$$

with probability at least $1 - d^{-10}$, we have for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$,

$$\Omega \left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \right) \cdot \mathbf{I} \preceq \nabla^2 f_n(\mathbf{W}) \preceq \Omega(1) \cdot \mathbf{I},$$

where $r_{\text{FCN}} := \frac{c_2}{\sqrt{K}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}$.

- For CNN, assume $\|\mathbf{w}^*\|_2 \leq 1$. There exist constants c_3 and c_4 such that as soon as sample size

$$n_{\text{CNN}} \geq c_3 \cdot dK^5 \log^2 d \cdot \left(\frac{1}{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)} \right)^2,$$

with probability at least $1 - d^{-10}$, we have for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})$,

$$\Omega \left(\frac{1}{K} \cdot \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2) \right) \cdot \mathbf{I} \preceq \nabla^2 f_n(\mathbf{w}) \preceq \Omega(K) \cdot \mathbf{I},$$

where $r_{\text{CNN}} := \frac{c_4}{K^2} \cdot \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)$.

We note that for FCN (3.1), all column permutations of \mathbf{W}^* are equivalent global minimum of the loss function, and Theorem 3 applies to all such permutation matrices of \mathbf{W}^* . The proof of Theorem 3 is outlined in Appendix B.2.

A pivot observation from the lower bound of the Hessian is that the sign of $\rho(\cdot)$ will determine whether the Hessian is positive definite or not, since K, κ, λ are all positive. We depict $\rho(\sigma)$ as a function of σ in a certain range for the sigmoid activation in Fig. 3.2. It can be seen from Fig. 3.2 that $\rho(\sigma)$ is monotonic increasing when σ increases, and we have $\rho(\sigma) > 0$ as long as $\sigma > 0$. When \mathbf{W}^* is orthogonal, κ and λ are both 1, $\rho(\sigma)$ is a constant, hence the lower bound of Hessian is on the order of $\frac{1}{K^2}$ for FCN. However, in the worst case where the columns of \mathbf{W}^* is linear dependent, then $\kappa, \lambda, \rho(\sigma)$ are infinite, and the local strong convexity doesn't hold for FCN case. Furthermore, the value of $\rho_{\text{CNN}}(\sigma)$ is much larger than $\rho_{\text{FCN}}(\sigma)$ for the same input.

Theorem 3 guarantees that for both FCN (3.1) and CNN (3.2) the Hessian of the empirical cross-entropy loss function $f_n(\mathbf{W})$ is positive definite in a neighborhood of the ground truth \mathbf{W}^* , as long as the sample size n is sufficiently large and the columns of \mathbf{W}^* are linearly independent. The bounds in Theorem 3 depend on the dimension parameters of the network (n and K), as well as the ground truth ($\rho_{\text{FCN}}(\sigma_K), \lambda, \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)$).

3.2.2 Performance Guarantees of GD

For the classification problem, due to the nature of quantized labels, \mathbf{W}^* is no longer a critical point of $f_n(\mathbf{W})$. By the strong convexity of the empirical risk function $f_n(\mathbf{W})$ in the local neighborhood of \mathbf{W}^* , there can exist at most one critical point in $\mathbb{B}(\mathbf{W}^*, r)$, which is the unique local minimizer in $\mathbb{B}(\mathbf{W}^*, r)$ if it exists. The following theorem shows that there indeed exists such a critical point $\widehat{\mathbf{W}}_n$, which is provably close to the ground truth \mathbf{W}^* , and gradient descent converges linearly to $\widehat{\mathbf{W}}_n$.

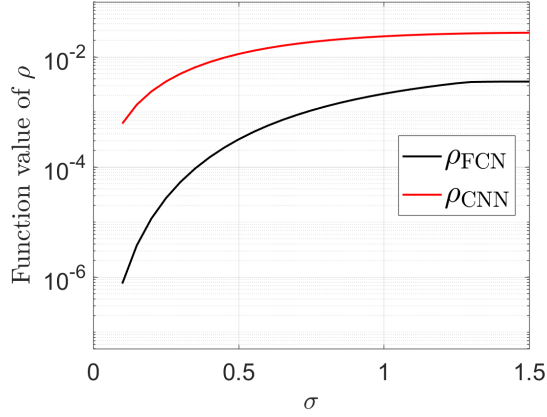


Figure 3.2: Illustration of $\rho(\sigma)$ for both FCN and CNN with the sigmoid activation.

Theorem 4 (Performance Guarantees of Gradient Descent). *Assume the assumptions in Theorem 3 hold. Under the event that local strong convexity holds,*

- for FCN, there exists a critical point in $\mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$ such that

$$\|\widehat{\mathbf{W}}_n - \mathbf{W}^*\|_{\text{F}} \leq c_1 \frac{K^{9/4} \kappa^2 \lambda}{\rho_{\text{FNN}}(\sigma_K)} \sqrt{\frac{d \log n}{n}},$$

and if the initial point $\mathbf{W}_0 \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$, GD converges linearly to $\widehat{\mathbf{W}}_n$, i.e.

$$\|\mathbf{W}_t - \widehat{\mathbf{W}}_n\|_{\text{F}} \leq \left(1 - \frac{c_2 \eta \rho_{\text{FCN}}(\sigma_K)}{K^2 \kappa^2 \lambda}\right)^t \|\mathbf{W}_0 - \widehat{\mathbf{W}}_n\|_{\text{F}},$$

for $\eta \leq c_3$, where c_1, c_2, c_3 are constants;

- for CNN, there exists a critical point in $\mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})$ such that

$$\|\widehat{\mathbf{w}}_n - \mathbf{w}^*\|_2 \leq c_4 \frac{K}{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)} \cdot \sqrt{\frac{d \log n}{n}},$$

and if the initial point $\mathbf{w}_0 \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})$, GD converges linearly to $\widehat{\mathbf{w}}_n$, i.e.

$$\|\mathbf{w}_t - \widehat{\mathbf{w}}_n\|_2 \leq \left(1 - \frac{c_5 \eta \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K}\right)^t \|\mathbf{w}_0 - \widehat{\mathbf{w}}_n\|_2,$$

for $\eta \leq c_6/K$, where c_4, c_5, c_6 are constants.

Similarly to Theorem 3, for FCN (3.1) Theorem 4 also holds for all column permutations of \mathbf{W}^* . The proof can be found in Appendix B.3. Theorem 4 guarantees that the existence of critical points in the local neighborhood of the ground truth, which GD converges to, and also shows that the critical points converge to the ground truth \mathbf{W}^* at the rate of $O(K^{9/4}\sqrt{d\log n/n})$ for FCN (3.1) and $O\left(K\sqrt{d\log n/n}\right)$ for CNN(3.2) with respect to increasing the sample size n . Therefore, \mathbf{W}^* can be recovered consistently as n goes to infinity. Moreover, for both FCN (3.1) and CNN (3.2) gradient descent converges linearly to $\widehat{\mathbf{W}}_n$ (or resp. $\widehat{\mathbf{w}}_n$) at a linear rate, as long as it is initialized in the basin of attraction. To achieve ϵ -accuracy, i.e. $\|\mathbf{W}_t - \widehat{\mathbf{W}}_n\|_{\text{F}} \leq \epsilon$ (or resp. $\|\mathbf{w}_t - \widehat{\mathbf{w}}_n\|_2 \leq \epsilon$), it requires a computational complexity of $O(ndK^4 \log(1/\epsilon))$ (or resp. $O(ndK^2 \log(1/\epsilon))$), which is linear in n , d and $\log(1/\epsilon)$.

3.3 Initialization via Tensor Method

Our initialization adopts the tensor method proposed in [7]. The initialization method works for the FCN model directly, and works for the CNN model with slight modification as presented in [79]. To avoid unnecessary repetitions from the previous work, we focus on the FCN case to outline the algorithm and remark the difference. We recommend the readers refer to [7, 79] for more details.

3.3.1 Preliminary and Algorithm

We start with introducing the necessary definitions which can be found in [7]. We first define a product $\tilde{\otimes}$ as follows. If $\mathbf{v} \in \mathbb{R}^d$ is a vector and \mathbf{I} is the identity matrix, then $\mathbf{v} \tilde{\otimes} \mathbf{I} = \sum_{j=1}^d [\mathbf{v} \otimes \mathbf{e}_j \otimes \mathbf{e}_j + \mathbf{e}_j \otimes \mathbf{v} \otimes \mathbf{e}_j + \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}]$. If \mathbf{M} is a symmetric

rank- r matrix factorized as $\mathbf{M} = \sum_{i=1}^r \mathbf{s}_i \mathbf{v}_i \mathbf{v}_i^\top$ and \mathbf{I} is the identity matrix, then

$$\mathbf{M} \tilde{\otimes} \mathbf{I} = \sum_{i=1}^r \mathbf{s}_i \sum_{j=1}^d \sum_{l=1}^6 \mathbf{A}_{l,i,j}, \quad (3.7)$$

where $\mathbf{A}_{1,i,j} = \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_j$, $\mathbf{A}_{2,i,j} = \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{e}_j$, $\mathbf{A}_{3,i,j} = \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{e}_j$, $\mathbf{A}_{4,i,j} = \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}_i$, $\mathbf{A}_{5,i,j} = \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{v}_i$ and $\mathbf{A}_{6,i,j} = \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{v}_i$.

This allows us to introduce the following quantities.

Definition 3. Define $\mathbf{M}_1, \mathbf{M}_2, \mathbf{M}_3, \mathbf{M}_4$ and $m_{1,i}, m_{2,i}, m_{3,i}, m_{4,i}$ as follows:

$$\mathbf{M}_1 = \mathbb{E}[y \cdot \mathbf{x}],$$

$$\mathbf{M}_2 = \mathbb{E}[y \cdot (\mathbf{x} \otimes \mathbf{x} - \mathbf{I})],$$

$$\mathbf{M}_3 = \mathbb{E}[y \cdot (\mathbf{x}^{\otimes 3} - \mathbf{x} \tilde{\otimes} \mathbf{I})],$$

$$\mathbf{M}_4 = \mathbb{E}[y \cdot (\mathbf{x}^{\otimes 4} - (\mathbf{x} \otimes \mathbf{x}) \tilde{\otimes} \mathbf{I} + \mathbf{I} \tilde{\otimes} \mathbf{I})],$$

$$m_{l,i} = g_{l,i}(\|\mathbf{w}_i^*\|), \forall l = 0, 1, 2, 3, 4,$$

where $g_{1,i}(\sigma) = \gamma_1(\sigma)$, $g_{2,i}(\sigma) = \gamma_2(\sigma) - \gamma_0(\sigma)$, $g_{3,i}(\sigma) = \gamma_3(\sigma) - 3\gamma_1(\sigma)$, $g_{4,i}(\sigma) = \gamma_4(\sigma) + 3\gamma_0(\sigma) - 6\gamma_2(\sigma)$, and $\gamma_j(\sigma) = \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\phi(\sigma \cdot z) z^j]$, $\forall j = 0, 1, 2, 3, 4$.

We further define a tensor operation as follows. For a tensor $\mathbf{T} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ and three matrices $\mathbf{A} \in \mathbb{R}^{n_1 \times d_1}$, $\mathbf{B} \in \mathbb{R}^{n_2 \times d_2}$, $\mathbf{C} \in \mathbb{R}^{n_3 \times d_3}$, the (i, j, k) -th entry of the tensor $\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C})$ is given by

$$\sum_{i'}^{n_1} \sum_{j'}^{n_2} \sum_{k'}^{n_3} \mathbf{T}_{i',j',k'} \mathbf{A}_{i',i} \mathbf{B}_{j',j} \mathbf{C}_{k',k}. \quad (3.8)$$

Armed this with definition, we define the following useful quantities.

Definition 4. Let $\boldsymbol{\alpha} \in \mathbb{R}^d$ denote a randomly picked vector. We define \mathbf{P}_2 and \mathbf{P}_3 as follows: $\mathbf{P}_2 = \mathbf{M}_{j_2}(\mathbf{I}, \mathbf{I}, \boldsymbol{\alpha}, \dots, \boldsymbol{\alpha})$, where $j_2 = \min\{j \geq 2 | \mathbf{M}_j \neq 0\}$, and $\mathbf{P}_3 = \mathbf{M}_{j_3}(\mathbf{I}, \mathbf{I}, \mathbf{I}, \boldsymbol{\alpha}, \dots, \boldsymbol{\alpha})$, where $j_3 = \min\{j \geq 3 | \mathbf{M}_j \neq 0\}$.

We further denote $\bar{\mathbf{w}} = \mathbf{w}/\|\mathbf{w}\|$. An important implication of Definition 3 and 4 is that the non-zero matrix \mathbf{P}_2 and non-zero tensor \mathbf{P}_3 is in the form of $\sum_{i=1}^K m_{j_2,i} (\alpha^\top \bar{\mathbf{w}}_i^*)^{j_2-2} \bar{\mathbf{w}}_i^{*\otimes 2}$, $\sum_{i=1}^K m_{j_3,i} (\alpha^\top \bar{\mathbf{w}}_i^*)^{j_3-3} \bar{\mathbf{w}}_i^{*\otimes 3}$, see [7, Claim 5.5]. The basic strategy is to extract the direction, magnitude information from the empirical version of \mathbf{P}_2 and \mathbf{P}_3 . Hence estimating \mathbf{W}^* can be decomposed as the following two steps.

- Estimate the direction of each column of \mathbf{W}^* by decomposing \mathbf{P}_2 to approximate the subspace spanned by $\{\bar{\mathbf{w}}_1^*, \bar{\mathbf{w}}_2^*, \dots, \bar{\mathbf{w}}_K^*\}$ (denoted by \mathbf{V}), then reduce the third-order tensor \mathbf{P}_3 to a lower-dimension tensor $\mathbf{R}_3 = \mathbf{P}_3(\mathbf{V}, \mathbf{V}, \mathbf{V}) \in \mathbb{R}^{K \times K \times K}$, and apply non-orthogonal tensor decomposition on \mathbf{R}_3 to output the estimate $s_i \mathbf{V}^\top \bar{\mathbf{w}}_i^*$, where $s_i \in \{1, -1\}$ is a random sign.
- Approximate the magnitude of \mathbf{w}_i^* and the sign s_i by solving a linear system of equations.

The initialization algorithm based on the tensor method is outlined in Algorithm 2. For more implementation details about Algorithm 2, e.g., power method, we refer to [7].

Algorithm 2 Initialization via Tensor Method

Require: Partition n pairs of data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ into three subsets $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3$.

Ensure:

- 1: Estimate $\hat{\mathbf{P}}_2$ of \mathbf{P}_2 from data set \mathcal{D}_1 .
 - 2: $\mathbf{V} \leftarrow \text{POWERMETHOD}(\hat{\mathbf{P}}_2, K)$.
 - 3: Estimate $\hat{\mathbf{R}}_3$ of $\mathbf{P}_3(\mathbf{V}, \mathbf{V}, \mathbf{V})$ from data set \mathcal{D}_2 .
 - 4: $\{\hat{\mathbf{u}}_i\}_{i \in [K]} \leftarrow \text{KCL}(\hat{\mathbf{R}}_3)$.
 - 5: $\{\mathbf{w}_i^{(0)}\}_{i \in [K]} \leftarrow \text{RECMAG}(\mathbf{V}, \{\hat{\mathbf{u}}_i\}_{i \in [K]}, \mathcal{D}_3)$.
-

3.3.2 Performance Guarantee of Initialization

For the classification problem, we make the following technical assumptions, similarly to [7, Assumption 5.3] for the regression problem.

Assumption 1. *The activation function $\phi(z)$ satisfies the following conditions:*

1. *If $M_j \neq 0$, then*

$$\sum_{i=1}^K m_{j,i} \left(\mathbf{w}_i^{\star\top} \boldsymbol{\alpha} \right)^{j-2} \overline{\mathbf{w}_i^{\star} \mathbf{w}_i^{\star\top}} \neq \mathbf{0},$$

$$\sum_{i=1}^K m_{j,i} \left(\overline{\mathbf{w}_i^{\star\top} \boldsymbol{\alpha}} \right)^{j-3} (\mathbf{V}^\top \overline{\mathbf{w}_i^{\star}}) \text{vec}((\mathbf{V}^\top \overline{\mathbf{w}_i^{\star}})(\mathbf{V}^\top \overline{\mathbf{w}_i^{\star}})^\top)^\top \neq \mathbf{0},$$

for $j \geq 3$.

2. *At least one of M_3 and M_4 is non-zero.*

Assumption 1 is to guarantee that the key terms still contain the magnitude information about \mathbf{w}_j^* . It can be verified that for sigmoid activation $m_{3,i}$ is non-zero for $\sigma > 0$, hence it will satisfy Assumption 1. Furthermore, we do not require the homogeneous assumption (i.e., $\phi(az) = a^p z$ for an integer p) required in [7], which can be restrictive. Instead, we assume the following condition on the curvature of the activation function around the ground truth, which holds for a larger class of activation functions such as sigmoid and tanh.

Assumption 2. *Let l_1 be the index of the first nonzero M_i where $i = 1, \dots, 4$. For the activation function $\phi(\cdot)$, there exists a positive constant δ such that $g_{l_1,i}(\cdot)$ is strictly monotone over the interval $(\|\mathbf{w}_i^*\| - \delta, \|\mathbf{w}_i^*\| + \delta)$, and the derivative of $g_{l_1,i}(\cdot)$ is lower bounded by some constant for all i .*

It can be numerically verified that sigmoid activation will also satisfy Assumption 2. We next present the performance guarantee for the initialization algorithm in the following theorem.

Theorem 5. *For the classification model (3.1), under Assumptions 1 and 2, for any $0 < \epsilon < 1$ and $\zeta > 1$, if the sample size $n \geq d \cdot \text{poly}(K, \kappa, \zeta, \log d, 1/\epsilon)$, then the output $\mathbf{W}_0 \in \mathbb{R}^{d \times K}$ of Algorithm 2 satisfies*

$$\|\mathbf{W}_0 - \mathbf{W}^*\|_{\text{F}} \leq \epsilon \text{poly}(K, \kappa) \|\mathbf{W}^*\|_{\text{F}}, \quad (3.9)$$

with probability at least $1 - d^{-\Omega(\zeta)}$.

The proof of Theorem 5 consists of (a) showing the estimation of the direction of \mathbf{W}^* is sufficiently accurate and (b) showing the approximation of the norm of \mathbf{W}^* is accurate enough. The proof of part (a) is the same as that in [7], but our argument in part (b) is different, where we relax the homogeneous assumption on activation functions. More details can be found in the supplementary materials in Appendix B.5.

3.4 Numerical Experiments

For FCN, we first implement gradient descent to verify that the empirical risk function is strongly convex in the local region around \mathbf{W}^* . If we initialize multiple times in such a local region, it is expected that gradient descent converges to the same critical point $\widehat{\mathbf{W}}_n$, with the same set of training samples. Given a set of training samples, we randomly initialize multiple times, and then calculate the variance of the output of gradient descent. Denote the output of the ℓ th run as $\hat{\mathbf{w}}_n^{(\ell)} = \text{vec}(\widehat{\mathbf{W}}_n^{(\ell)})$ and the mean of the runs as $\bar{\mathbf{w}}$. The error is calculated as $\text{SD}_n = \sqrt{\frac{1}{L} \sum_{\ell=1}^L \|\hat{\mathbf{w}}_n^{(\ell)} - \bar{\mathbf{w}}\|^2}$, where $L = 20$ is the total number of random initializations. Adopted from [50], it

quantifies the standard deviation of the estimator $\widehat{\mathbf{W}}_n$ under different initializations with the same set of training samples. We say an experiment is successful, if $\text{SD}_n \leq 10^{-4}$. We generate the ground truth \mathbf{W}^* from Gaussian matrices, and the training samples are generated using the FCN (3.1). Figure 3.3 (a) shows the successful rate of gradient descent by averaging over 50 sets of training samples for each pair of n and d , where $K = 3$ and $d = 15, 20, 25$ respectively. The maximum iterations for gradient descent is set as $\text{iter}_{\max} = 3500$. It can be seen that as long as the sample complexity is large enough, gradient descent converges to the same local minima with high probability.

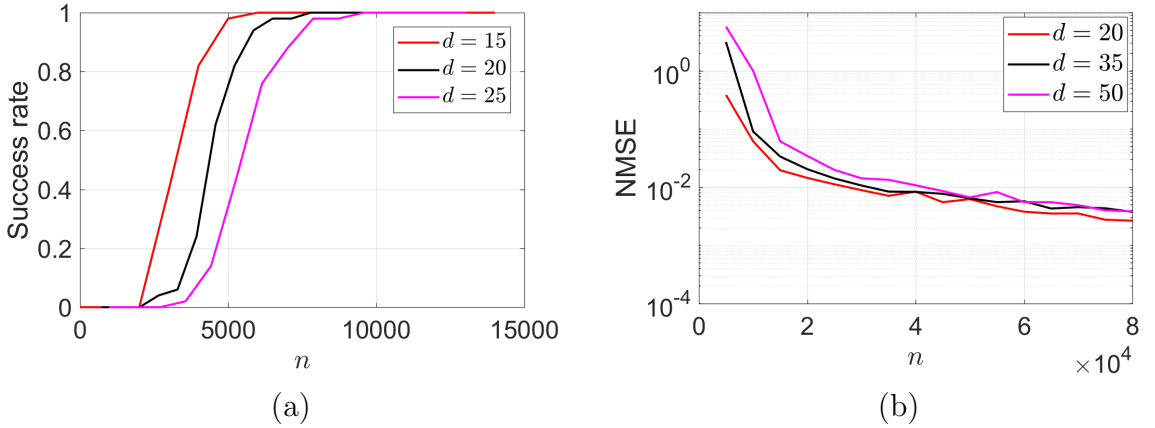


Figure 3.3: For FCN (3.1) fix $K = 3$. (a) Success rate of converging to the same local minima with respect to the sample complexity for various d with threshold 10^{-4} ; (b) Average estimation error of gradient descent in a local neighborhood of the ground truth with respect to the sample complexity for various d . The x-axis is scaled to illuminate the correct scaling between n and d .

We next show that the statistical accuracy of the local minimizer for gradient descent if it is initialized close enough to the ground truth. Suppose we initialize around

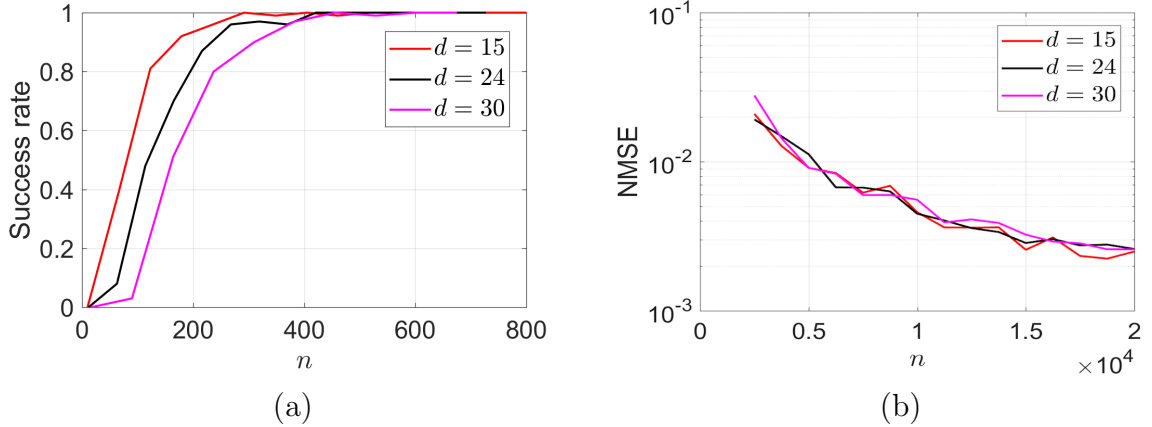


Figure 3.4: For CNN (3.2), fix $K = 3$. (a) Success rate of converging to the same local minima with respect to the sample complexity for various d with threshold 10^{-14} ; (b) Average estimation error of gradient descent in a local neighborhood of the ground truth with respect to the sample complexity for various d . The x-axis is scaled to illuminate the correct scaling between n and d .

the ground truth such that $\|\mathbf{W}_0 - \mathbf{W}^*\|_F \leq 0.1 \cdot \|\mathbf{W}^*\|_F$. We calculate the average estimation error as $\sum_{\ell=1}^L \|\widehat{\mathbf{W}}_n^{(\ell)} - \mathbf{W}^*\|_F^2 / (L \|\mathbf{W}^*\|_F^2)$ over $L = 100$ Monte Carlo simulations with random initializations. Fig. 3.3 (b) shows the average estimation error with respect to the sample complexity when $K = 3$ and $d = 20, 35, 50$ respectively. It can be seen that the estimation error decreases gracefully as we increase the sample size and matches with the theoretical prediction of error rates reasonably well.

Similarly, for CNN, we first verify that the empirical risk function is locally strongly convex using the same method as before. We generate the entries of true weights \mathbf{w}^* from standard Gaussian distribution, and generate the training samples using the CNN model (3.2). In Fig. 3.4 (a), we say an experiment is successful if $SD_n \leq 10^{-14}$, and the successful rate is calculated over 100 sets of training samples

with $K = 3$ and $d = 15, 24, 30$ respectively. Then we verify the performance of gradient descent in Fig. 3.4 (b). Suppose we initialized in the neighborhood of \mathbf{w}^* , i.e., $\|\mathbf{w}_0 - \mathbf{w}^*\|_2 \leq 0.9 \cdot \|\mathbf{w}^*\|_2$, for fixed d, K, n , the average error is calculated over $L = 100$ Monte Carlo simulations. It can be seen that the error decreases as we increase the number of samples.

Chapter 4: Guaranteed Recovery of CNN with ReLU Activations

In this chapter, we study the model recovery problem when the data is generated by a one-hidden-layer non-overlap convolutional neural network with ReLU activations, and the goal is to recover the weights of the neural network. We first characterize the landscape of the population risk function $L(\mathbf{w})$ by exploiting the second-order property of $L(\mathbf{w})$, and then we show that gradient descent with random initialization converges to the global minimum \mathbf{w}^* at a much faster rate than previously established. Furthermore, we also study the empirical risk function which is practically used in training a neural network. We show that with a well designed initialization, gradient descent converges linearly to the true weights \mathbf{w}^* with high probability.

4.1 Problem Formulation

Given an input $\mathbf{x} \in \mathbb{R}^d$, the output of a one-hidden-layer non-overlapping convolutional neural network as illustrated in Fig. 4.1 is given by

$$f(\mathbf{x}; \mathbf{w}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}^\top \mathbf{x}^{(k)}), \quad (4.1)$$

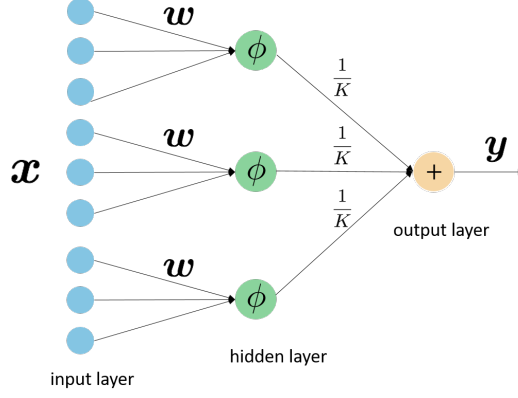


Figure 4.1: Illustration of a one-hidden-layer convolutional neural network without overlap

where K is the number of neurons, $\phi(\cdot)$ is the ReLU activation function, i.e., $\phi(z) = \max(0, z)$, $\mathbf{w} \in \mathbb{R}^m$ is the weight of the neural network, and the k th stride of \mathbf{x} be given as $\mathbf{x}^{(k)} = [x_{m(k-1)+1}, \dots, x_{m \cdot k}]^\top \in \mathbb{R}^m$. Here, we assume $d = m \cdot K$ for simplicity.

Suppose we are given a set of training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where y_i is generated by the neural network with \mathbf{w}^* , i.e., $y_i = f(\mathbf{x}_i; \mathbf{w}^*)$. Our goal is to learn the true neural network \mathbf{w}^* via risk minimization. Under the regression setting, the squared loss is generally adopted as the risk function. Two types of risk functions have been considered in the literature, i.e., the population risk function

$$L(\mathbf{w}) = \mathbb{E}_{\mathcal{D}} [(y - f(\mathbf{x}; \mathbf{w}))^2], \quad (4.2)$$

where \mathcal{D} denotes the joint distribution of (\mathbf{x}, y) , and the empirical risk function

$$L_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \mathbf{w}))^2. \quad (4.3)$$

In this work, we first consider minimizing the population risk function (4.2) via gradient descent, and then we further analyze the performance of gradient descent on minimizing the empirical risk function (4.3).

As shown in [71], even with infinitely many samples, the problem of minimizing the population risk without any assumption on the distribution \mathcal{D} is NP-Complete. Hence in this work we adopt the standard assumption as in [71, 73] that the input $\{\mathbf{x}_i\}_{i=1}^n$ are composed of i.i.d. standard Gaussian vectors.

4.2 Minimizing the Population Risk

To minimize (4.2), we implement the gradient descent (GD) algorithm which is an iterative method and update the parameter as follows.

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t^{(1)} \nabla L(\mathbf{w}_t), \quad (4.4)$$

where $\eta_t^{(1)}$ is the step size or the learning rate, and the subscript t denotes the t -th update. This algorithm and its variant (such as stochastic gradient descent) have been widely used in practical machine learning and deep learning.

The population risk function has been studied in [71], in which only the first-order property was exploited. We take another path by leveraging the second-order property of the population risk function to provide a refined analysis. According to [71], under the Gaussian input assumption, we calculate the closed form of the population risk function (4.2) (up to the additive factors in \mathbf{w}^*) as

$$L(\mathbf{w}) = \frac{1}{K^2} \left[\left(\frac{K^2 - K}{2\pi} + \frac{K}{2} \right) \|\mathbf{w}\|_2^2 - \frac{K^2 - K}{\pi} \|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2 - \frac{K}{\pi} \|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2 (\sin(\theta) + (\pi - \theta) \cos(\theta)) \right], \quad (4.5)$$

where $\theta = \arccos\left(\frac{\mathbf{w}^\top \mathbf{w}^*}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2}\right)$ denotes the angle between the vectors \mathbf{w} and \mathbf{w}^* .

The population risk function (4.5) in general is nonconvex, and with the closed form expression of the population risk function, the critical points can be characterized. According to [71, Lemma 5.1], the population risk function (4.5) has three critical points, i.e.,

- (a) A local maximum at $\mathbf{w} = \mathbf{0}$;
- (b) A unique global minimum at $\mathbf{w} = \mathbf{w}^*$;
- (c) A degenerate saddle point at $\mathbf{w} = -\frac{K^2-K}{K^2+(\pi-1)K}\mathbf{w}^*$.

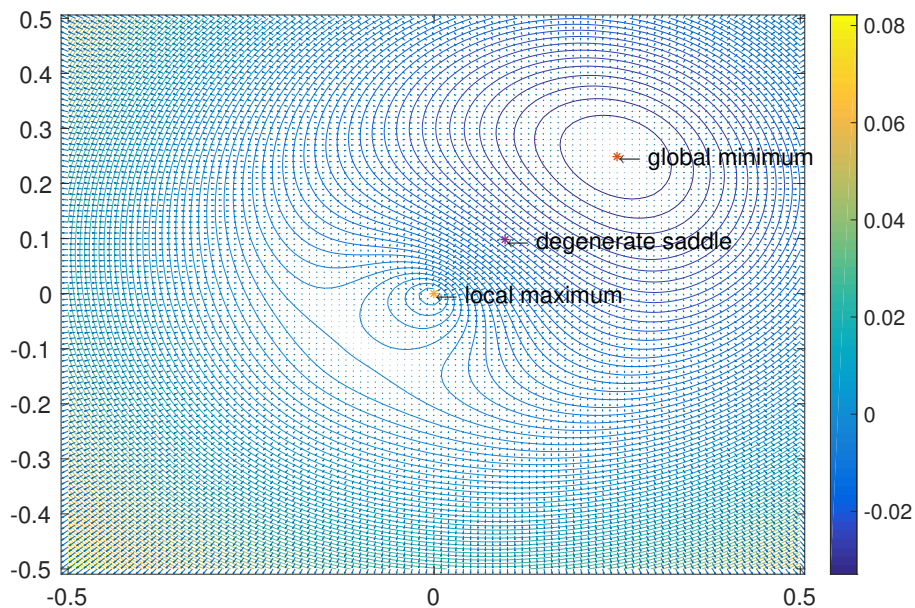


Figure 4.2: Contour plot of the population risk function (4.5) when $m = 2$ and $K = 3$.

We illustrate the landscape of the population risk function (4.5) in Fig 4.2 and mark the three critical points. With the aid of such critical point characterization, [71]

further showed that gradient descent with random initialization converges to a point with ϵ -accuracy to \mathbf{w}^* after $O\left(\frac{1}{\epsilon^4}\right)$ steps for which we suspect is not optimal. Moreover, the analysis of the critical points is not sufficient to fully understand the optimization landscape. In particular, we are interested in understanding whether the population risk is strongly convex in a region around the ground truth \mathbf{w}^* ; as demonstrated in Chapter 3 that the nonconvex objective functions for learning one-hidden-layer neural network usually have benign geometric properties locally. Hence, beyond the analysis of the critical points which only requires the first-order information of the population risk, we further analyze the Hessian of the population risk. We show that the Hessian of (4.2) can be lower and upper bounded in the local neighborhood of \mathbf{w}^* with a radius $\Omega\left(\frac{1}{K}\right)$. With the aid of this property, we are able to improve the performance analysis of gradient descent and show that with random initialization, gradient descent converges much faster to \mathbf{w}^* .

Before presenting the result, we denote $\mathbb{B}(\mathbf{w}^*, r)$ as a Euclidean ball centered at $\mathbf{w}^* \in \mathbb{R}^m$ with a radius r , i.e.,

$$\mathbb{B}(\mathbf{w}^*, r) = \{\mathbf{w} \in \mathbb{R}^m : \|\mathbf{w} - \mathbf{w}^*\|_2 \leq r\}. \quad (4.6)$$

The second-order property of the population risk function (4.2) is summarized as follows.

Theorem 6. *Consider the regression model with the ReLU activation function (4.1), and the population risk (4.2). If the input follows a Gaussian distribution, i.e., $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then the following inequality*

$$\frac{1}{4K} \cdot \mathbf{I} \preceq \nabla^2 L(\mathbf{w}) \preceq 3 \cdot \mathbf{I} \quad (4.7)$$

holds for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)$, as long as $\mathbf{w} \neq -\lambda \cdot \mathbf{w}^*$, where λ is any nonnegative constant and $r := \frac{1}{3K+3}$.

The proof of Theorem 6 is provided in Appendix C.2. Theorem 6 guarantees that the population loss is strongly convex in the neighborhood of \mathbf{w}^* , and the size of the neighborhood is roughly $O\left(\frac{1}{K}\right)$ which shrinks when the number of neurons increases. The implication of such a local strong convexity is that once initialized in the local region, if the updates of gradient descent \mathbf{w}_t never lies on the line of $-\mathbf{w}^*$, gradient descent converges to the unique critical point \mathbf{w}^* at a linear rate. In fact, it has been shown in [71] that the angle between the updates \mathbf{w}_t and \mathbf{w}^* is decreasing, and thus \mathbf{w}_t never lies on the line of $-\mathbf{w}^*$.

It has been shown in [71] that the estimation error $\|\mathbf{w}_t - \mathbf{w}^*\|_2$ keeps decreasing during the execution of gradient descent. The analysis can be used to guarantee that \mathbf{w}_t enters such a local region. Hence, combining these two phases yields the following guarantee for (4.4).

Theorem 7. *Assume $\|\mathbf{w}^*\|_2 = 1$, and suppose the initial point \mathbf{w}_0 is uniformly drawn from the unit sphere. Further set the learning rate $\eta_t^{(1)} = \eta < \frac{1}{3}$. Then the gradient descent updates \mathbf{w}_t satisfy*

$$\|\mathbf{w}_t - \mathbf{w}^*\|_2 \leq \epsilon \tag{4.8}$$

after $O\left(K^4 + K \cdot \log\left(\frac{1}{\epsilon}\right)\right)$ iterations.

Theorem 7 shows that even though the population loss (4.5) is nonconvex, gradient descent still converges fast to the global minimum with random initialization. The theoretical guarantee suggests that this convergence exhibits two phases: first the estimation error contracts relatively slowly; once it enters the local neighborhood

of \mathbf{w}^* after $O(K^4)$ steps, the convergence rate speeds up. Comparing to the result of [71] which requires $O\left(\frac{1}{\epsilon^4}\right)$ iterations to guarantee the ϵ -accuracy, our result shows that gradient descent converges much faster. This improvement highly depends on the property of the Hessian that we characterize. We sketch the proof of Theorem 7 as follows.

Proof. Let \mathbf{w}_t be the estimator at the t -th iteration. Due to the proof of [71, Theorem 5.2], we have that with constant step size $0 < \eta < 1$,

$$\|\mathbf{w}_t - \mathbf{w}^*\|_2 \leq O\left(\frac{1}{K}\right) \quad (4.9)$$

holds after $O(K^4)$ iterations, and $\frac{\mathbf{w}^\top \mathbf{w}^*}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2} \neq -1$. Hence, after $O(K^4)$ iterations, w_t enters the local region where

$$\frac{1}{4K} \cdot \mathbf{I} \preceq \nabla^2 L(\mathbf{w}) \preceq 3 \cdot \mathbf{I} \quad (4.10)$$

holds according to Theorem 6. Then following the same proof as in [59, Lemma 1], we obtain the following error contraction result,

$$\|\mathbf{w}_{t+1} - \mathbf{w}^*\|_2 \leq \left(1 - \frac{1}{K}\right) \|\mathbf{w}_t - \mathbf{w}^*\|_2, \quad (4.11)$$

as long as $\eta < \frac{1}{3}$. Hence after $O(K^4)$ iterations, gradient descent converges to \mathbf{w}^* linearly. \square

4.3 Minimizing the Empirical Risk

In practice, it is common to minimize the empirical risk function formed with finite samples. Hence, in this section we study the empirical risk function. Particularly, we implement the same gradient descent algorithm to minimize (4.3), which updates

\mathbf{w}_{t+1} via

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t^{(2)} \nabla L_n(\mathbf{w}_t), \quad (4.12)$$

where $\eta_t^{(2)}$ is the step size or the learning rate. In this section, we analyze the performance of gradient descent (4.12). Note that throughout the execution of GD, the same set of training samples is used which is the standard implementation of gradient descent. Consequently the analysis is challenging due to the statistical dependence of the iterates with the data. Another challenge of analyzing the performance of GD on minimizing the empirical risk function is due to the non-smoothness of the empirical risk function because ReLU is not differentiable at 0. In our algorithm we define the gradient of ReLU activation $\phi(\cdot)$ as $\phi'(\mathbf{x}) = 1_{\{x>0\}}$. Then we calculate the gradient of the empirical risk function (4.3) as

$$\nabla L_n(\mathbf{w}) = 2 \cdot \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \left(\phi(\mathbf{w}^\top \mathbf{x}_i^{(j)}) - \phi(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}) \right) \phi'(\mathbf{w}^\top \mathbf{x}_i^{(l)}) \mathbf{x}_i^{(l)}. \quad (4.13)$$

Since the empirical risk function (4.3) is highly non-convex and non-smoothness, it is not possible to exploit the geometric property of (4.3) by analyzing its Hessian directly. As the analysis developed in Chapter 3, the empirical risk function usually concentrates around the population risk function when the sample size is large enough. However, due to the unboundedness of ReLU activation, it is very difficult to implement this approach and guarantee the concentration.

Instead, we resort to a so-called regularity condition, stated as follows, which is widely used to establish the geometric property of the loss function.

Definition 5 (Regularity condition [59]). *The function $f(\cdot)$ is said to obey the regularity condition $RC(\mu, \lambda, \zeta)$ for some $\mu, \lambda, \zeta > 0$ if*

$$2 \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\nabla f(\mathbf{w})\|_2^2 + \lambda \|\mathbf{w} - \mathbf{w}^*\|_2^2 \quad (4.14)$$

for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, \zeta)$.

Note that this condition does not require $\nabla f(\cdot)$ to be differentiable and it does not require the loss function $f(\cdot)$ to be convex, which is applicable to our case. Moreover, such a condition implies that at any point $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, \zeta)$ the associated gradient $\nabla f(\cdot)$ is positively correlated with the estimation error $\mathbf{w} - \mathbf{w}^*$, and hence the update rule (4.12) reduces the error $\mathbf{w} - \mathbf{w}^*$. Next we show that the regularity condition indeed holds for (4.13) in Lemma 1. By leveraging such a geometric property, we show that gradient descent (4.12) converges linearly.

Lemma 1. *For all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)$ with $r := \frac{1}{K^{\frac{3}{2}}}\|\mathbf{w}^*\|_2$, the following regularity condition*

$$\langle \nabla L_n(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \frac{1}{72K^2} \|\nabla L_n(\mathbf{w})\|_2^2 + \frac{1}{4K^2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \quad (4.15)$$

holds with probability at least $1 - \frac{1}{d^{10}}$, as long as the sample complexity satisfies $n \geq c \cdot mK^2 \cdot \log(n)$ for some sufficiently large constant c .

The proof of Lemma 1 can be found in Appendix C.3.1, in which we do not optimize the constant. Lemma 1 captures the particular geometric property of the empirical risk function within a local neighborhood of \mathbf{w}^* . Utilizing such a geometric property, we next show the performance guarantee of gradient descent (4.12) as below.

Theorem 8. *Consider the problem of minimizing (4.3) with gradient descent (4.12). Assume the input $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the initial point satisfies $\mathbf{w}_0 \in \mathbb{B}(\mathbf{w}^*, r)$ with*

$r := \frac{1}{K^{\frac{3}{2}}}\|\mathbf{w}^*\|_2$. Then, with the step size $\eta_t^{(2)} = \frac{1}{4K^2}$, the following inequality

$$\|\mathbf{w}_t - \mathbf{w}^*\|_2^2 \leq \left(1 - \frac{1}{288K^4}\right)^t \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 \quad (4.16)$$

holds with probability at least $1 - \frac{1}{d^{10}}$ as long as the sample complexity satisfies $n \geq c \cdot mK^2 \cdot \log(n)$ for some sufficiently large constant c .

Theorem 8 implies that in the finite sample regime, gradient descent with a well designed initial point still converges to the ground truth although the empirical risk function is not convex. Note that the performance guarantee requires a local initialization with radius $\frac{1}{K^{\frac{3}{2}}}\|\mathbf{w}^*\|_2$, the sample complexity is roughly on the order of $O(dK)$ where $d = mK$, and the sample complexity is linear in terms of both the input dimension d and the number of neurons K .

Proof. Following the same proof of [59, Lemma 2], setting $\eta_t^{(2)} = \frac{1}{4K^2}$, and then applying Lemma 1, we obtain

$$\|\mathbf{w}_{t+1} - \mathbf{w}^*\|_2^2 \leq \left(1 - \frac{1}{288K^4}\right) \|\mathbf{w}_t - \mathbf{w}^*\|_2^2. \quad (4.17)$$

□

4.4 Numerical Experiments

In this section, we implement gradient descent on synthetic data to demonstrate the statistical accuracy of the optimizer if \mathbf{w} is initialized close enough to the ground truth \mathbf{w}^* , i.e., $\|\mathbf{w}_0 - \mathbf{w}^*\|_2 \leq \frac{1}{K^{\frac{3}{2}}}\|\mathbf{w}^*\|_2$. We first generate $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$ as Gaussian random vector, i.e., each entry of \mathbf{x}_i follows a standard Gaussian distribution, and then we generate the output y_i via (4.1). Given the set of training samples, we randomly initialize \mathbf{w}_0 and calculate the error of the outputs of gradient descent at

each step. Denoting the output at the t -th step as \mathbf{w}_t , we define the normalized mean squared error (NMSE) as $\frac{\|\mathbf{w}_t - \mathbf{w}^*\|_2^2}{\|\mathbf{w}^*\|_2^2}$. We first initialize \mathbf{w}_0 as $\mathbf{w}_0 = \mathbf{w}^* + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \frac{1}{K^3} \|\mathbf{w}^*\|_2^2)$, such that $\|\mathbf{w}_0 - \mathbf{w}^*\|_2 \leq \frac{1}{K^{\frac{3}{2}}} \|\mathbf{w}^*\|_2$, and Fig 4.3 (a) shows that with local initialization, gradient descent converges linearly in log scale, and \mathbf{w}^* can be recovered exactly after the gradient descent converges.

Although the theoretical guarantee of gradient descent needs a local initialization, we observe an interesting phenomenon that gradient descent converges to the ground truth \mathbf{w}^* even without local initialization. We use the same set of training samples and run gradient descent starting from \mathbf{w}_0 which follows a Gaussian distribution. As can be seen from Fig 4.3 (b), gradient descent still converges to \mathbf{w}^* at a linear rate.

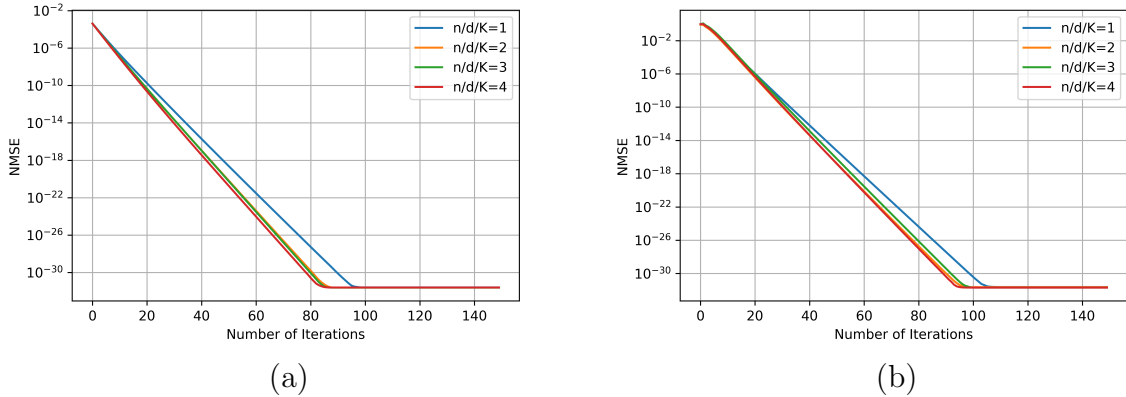


Figure 4.3: For CNN, fix $K = 10$, $m = 35$ and $d = 350$, the NMSE with respect to the number of steps of gradient descent for various n . (a) local initialization; (b) Gaussian initialization.

Chapter 5: Conclusion and Future Work

5.1 Concluding Remarks

In the study of line spectrum estimation from quantized measurement, we examined the effect of (heavy) quantization in spectral compressed sensing that is useful for understanding wideband spectral signal acquisition and processing. Our contributions are two-fold. We first derived the Cramér-Rao bound for parameter estimation with multiple complex sinusoids using quantized compressed linear measurements. This bound is instrumental in describing the trade-offs between bit depth and sample complexity at different SNR regimes. Such an estimation-theoretical perspective is independent of the algorithm and hasn't been exploited in the previous literature. Secondly, we developed algorithms for spectral-sparse signal recovery using quantized measurements via atomic norm minimization, which do not require knowledge of the quantizer in recovery. Under a mild separation condition, we establish that we can accurately recover a spectrally-sparse signal from the signs of $O(K \log n)$ random linear measurements. The proposed algorithm also can be extended to handle multiple signal snapshots. This generalizes the literature on one-bit compressed sensing to the

important class of spectrally sparse signals using atomic norms, and we carefully examined the performance of the proposed algorithms via numerical experiments. The content of this Chapter is published in [80, 81].

In the study of model recovery of neural-network-based models we recover an one-hidden-layer neural network using the cross-entropy loss in a multi-neuron classification problem. In particular, we have characterized the sample complexity to guarantee local strong convexity in a neighborhood (whose size we have characterized as well) of the ground truth when the training data are generated from a classification model for two types of neural network models: fully-connected network and non-overlapping convolutional network. This guarantees that with high probability, gradient descent converges linearly to the ground truth if initialized properly. The content of this chapter is summarized in a research paper and can be found in [82].

In the study of recovering a one-hidden-layer non-overlap convolutional neural network under regression setting, where the activation is the commonly used ReLU activation in practice. We first analyze the landscape of the population risk function more carefully and find that its Hessian is lower and upper bounded by some positive quantities. Leveraging such a good property, we improve the convergence rate of gradient descent in the existing literature significantly. Secondly, we analyze the non-smooth empirical risk function. We show that a so-called regularity condition holds uniformly in a neighborhood of \mathbf{w}^* with high probability as soon as the sample size is $O(mK^2)$. This further implies that gradient descent finds the global optimal \mathbf{w}^* at a linear convergence rate provided a good initialization.

5.2 Future Work

Along the study of estimating model parameters from coarse and nonlinear data, we obtained some insightful results and also found some interesting open problems. Hence we discuss several potential directions of future work in this section.

Alternative Algorithms for Spectrally-Sparse Signal Estimation

In the work of estimating spectrally sparse signal from its quantized linear measurements as described in Chapter 2, an alternative convex relaxation for spectrally-sparse signal recovery is based on Hankel matrix enhancement and nuclear norm minimization [83,84]. In the single vector case, instead of imposing the atomic norm regularizer as in (2.22), one may consider

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{C}^n} \frac{1}{2} \|\mathbf{x} - \mathbf{s}\|_2^2 + \tau_H \|\mathcal{H}(\mathbf{x})\|_*. \quad (5.1)$$

Here, $\mathcal{H}(\mathbf{x})$ denotes a Hankel matrix given as

$$\mathcal{H}(\mathbf{x}) = \begin{bmatrix} x_1 & x_2 & & & \\ x_2 & & \ddots & & \\ \vdots & & \ddots & & \\ x_{n_1} & x_{n_1+1} & \cdots & x_n \end{bmatrix},$$

where n_1 is set as $\lfloor n/2 \rfloor$ to make the matrix $\mathcal{H}(\mathbf{x})$ as square as possible, $\|\cdot\|_*$ is the nuclear norm, and τ_H is a regularization parameter. The preliminary numerical simulations suggest this method is also effective for promoting spectral sparsity, hence it's interesting to analyze (5.1) in the future work. Further more, since the Cramér-Rao bounds assume perfect knowledge of the quantizers, they may not be indicative to benchmark the performance of the atomic norm minimization algorithms as proposed in this paper, since these algorithms do not make use of such knowledge. In the future,

it might be interesting to develop estimation-theoretical bounds that only assume partial or little knowledge about the quantizer.

Learning One-Hidden-Layer Neural Networks with ReLU Activations

In the work of learning an one-hidden-layer neural network as discussed in Chapter 3 and 4, we considered different activation functions and network structures. However, ReLU activation is only considered with CNN network, an interesting future direction would be considering learn a FCN network with ReLU activation. Next we briefly introduce the problem. Suppose we are given a set of training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^d$ is drawn from an i.i.d. distribution, i.e., $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. And y_i is generated according to

$$y_i = \sum_{j=1}^K \phi(\mathbf{w}_j^{\star\top} \mathbf{x}_i), \quad (5.2)$$

where $\phi(x) = \max(0, x)$ is the ReLU activation, K is the number of neurons and we are interested in recovering $\mathbf{W}^* = [\mathbf{w}_1^*, \dots, \mathbf{w}_K^*] \in \mathbb{R}^{d \times K}$ by minimizing the following empirical risk function, i.e.

$$f_n(\mathbf{W}) = \frac{1}{2n} \sum_{i=1}^n \left(\sum_{j=1}^K \phi(\mathbf{w}_j^\top \mathbf{x}_i) - y_i \right)^2. \quad (5.3)$$

One challenge of solving this problem is that the loss function is non-convex and non-smoothness. It's hard to make use of the second order information about the empirical risk function. Previously, [69] studies the same problem, i.e. for the particular FCN network, they study minimizing the empirical risk function (5.3) via gradient descent. In [69] the authors showed that \mathbf{W}^* can be recovered up to certain estimation error. Such an estimation error is upper bounded by two quantities. The first quantity comes from the optimization aspect, it decreases extremely fast as the algorithm iterates

and it's not avoidable. However, the second term doesn't decrease as the algorithm iterates it diminishes only when the sample size increases to infinity. Hence in the finite sample regime gradient descent is not guaranteed to converge to \mathbf{W}^* exactly. Thus it's desired to further understand the performance of gradient descent, develop new analysis to avoid the second type of error and see if we can get an exact recover of \mathbf{W}^* in finite sample regime.

Another open question about learning the one-hidden-layer CNN with ReLU activations is how gradient descent performs with random initialization. We have already numerically demonstrate in Chapter 4 that gradient descent with random initialization still converges linearly to the ground truth. Moreover, it was shown in [71] that gradient descent with random initialization finds the global optimal of the population risk function. Hence, it's natural to postulate that apply gradient descent to learn CNN doesn't require well designed initialization.

Appendix A: Proofs for Chapter 2

A.1 Proof of Theorem 2

An alternative way to represent the atomic decomposition is to write it as an integration of certain point measure [28]. Define the representing measure of \mathbf{x}^* as

$$\mu(f) = \sum_{k=1}^K c_k \delta(f - f_k),$$

where $\delta(\cdot)$ is the delta function. Then we can rewrite \mathbf{x}^* as

$$\mathbf{x}^* = \int_0^1 \mathbf{v}(f) d\mu(f) = \sum_{k=1}^K c_k \mathbf{v}(f_k). \quad (\text{A.1})$$

Correspondingly, denote $\hat{\mu}(f)$ as the representing measure for the solution $\hat{\mathbf{x}}$ of (2.22), which means $\hat{\mathbf{x}} = \int_0^1 \mathbf{v}(f) d\hat{\mu}(f)$.

Denote the reconstruction error as $\mathbf{e} = \lambda \mathbf{x}^* - \hat{\mathbf{x}}$, and its representing measure is $\gamma = \lambda \mu - \hat{\mu}$. With these definitions, applying [28, Lemma 1], we can bound the error as [28]

$$\|\mathbf{e}\|_2^2 \leq \|\mathbf{e}\|_{\mathcal{A}}^* \left(\int_F |\gamma|(df) + I_0 + I_1 + I_2 \right), \quad (\text{A.2})$$

where $I_\ell = \sum_{k=1}^K I_\ell^k$, for $\ell = 0, 1, 2$, with $I_0^k = \left| \int_{N_k} \gamma(df) \right|$, $I_1^k = n \left| \int_{N_k} (f - f_k) \gamma(df) \right|$, $I_2^k = \frac{n^2}{2} \int_{N_k} (f - f_k)^2 |\gamma|(df)$, where $N_k = \{f \in \mathbb{T} : d(f, f_k) \leq 0.16/n\}$ as the neighborhoods around each frequency, and $F = \mathbb{T} \setminus \bigcap_{k=1}^K N_k$.

To bound the first term in (A.2), let us denote the deviation

$$\mathbf{w} = \mathbf{s} - \mathbb{E}[\mathbf{s}] = \mathbf{s} - \lambda \mathbf{x}^*, \quad (\text{A.3})$$

where $\mathbb{E}[\mathbf{w}] = 0$. We have

$$\begin{aligned} \|\mathbf{e}\|_{\mathcal{A}}^* &\leq \|\mathbf{w}\|_{\mathcal{A}}^* + \|\mathbf{s} - \hat{\mathbf{x}}\|_{\mathcal{A}}^* \\ &\leq \|\mathbf{w}\|_{\mathcal{A}}^* + \tau, \end{aligned} \quad (\text{A.4})$$

where the first line follows from the triangle inequality, and the second line follows from the optimality condition of the AST algorithm in (2.22) in the following lemma.

Lemma 2 (Optimality conditions [29]). *$\hat{\mathbf{x}}$ is the solution of (2.22) if and only if $\|\mathbf{s} - \hat{\mathbf{x}}\|_{\mathcal{A}}^* \leq \tau$, and $\langle \mathbf{s} - \hat{\mathbf{x}}, \hat{\mathbf{x}} \rangle = \tau \|\hat{\mathbf{x}}\|_{\mathcal{A}}$.*

Therefore, if we set $\tau \geq \eta \|\mathbf{w}\|_{\mathcal{A}}^*$, where $\eta \geq 1$ is some constant, then plugging this into (A.4) we can show that

$$\|\mathbf{e}\|_{\mathcal{A}}^* \leq (\eta^{-1} + 1)\tau \leq 2\tau. \quad (\text{A.5})$$

The second term in (A.2) can be bounded in exactly the same manner as in [28], as long as (A.5) holds. In effect, [28] proved the following bound, under the separation condition, with high probability we have

$$\left(\int_F |\gamma| (df) + I_0 + I_1 + I_2 \right) \leq C \frac{K\tau}{n}. \quad (\text{A.6})$$

The following lemma bounds $\|\mathbf{w}\|_{\mathcal{A}}^*$, whose proof is provided in Appendix A.2.

Lemma 3. *With probability at least $1 - 1/(\pi n \log n)$, we have*

$$\|\mathbf{w}\|_{\mathcal{A}}^* \leq C \cdot \sqrt{\frac{n \log n}{m}},$$

where C is some universal constant.

Therefore, set $\tau = C\eta\sqrt{n \log n/m}$, and plug (A.5) and (A.6) into (A.2), we have

$$\|\mathbf{e}\|_2^2 \leq C' \cdot \frac{K\tau^2}{n} \leq C' \frac{K \log n}{m}. \quad (\text{A.7})$$

which is equivalent to

$$\left\| \frac{\hat{\mathbf{x}}}{\lambda} - \mathbf{x}^* \right\|_2 \lesssim \frac{1}{\lambda} \sqrt{\frac{K \log n}{m}}.$$

The proof is complete.

A.2 Proof of Lemma 3

By definition, we can write $\|\mathbf{w}\|_{\mathcal{A}}^*$ as

$$\begin{aligned} \|\mathbf{w}\|_{\mathcal{A}}^* &= \sup_{f \in [0,1]} |\langle \mathbf{s} - \lambda \mathbf{x}^*, \mathbf{v}(f) \rangle| \\ &= \sup_{f \in [0,1]} |\langle \mathbf{s}, \mathbf{v}(f) \rangle - \mathbb{E}[\langle \mathbf{s}, \mathbf{v}(f) \rangle]| \\ &= \sup_{f \in [0,1]} |g_{\mathbf{x}^*}(f) - \mathbb{E}[g_{\mathbf{x}^*}(f)]| \end{aligned} \quad (\text{A.8})$$

where $g_{\mathbf{x}^*}(f) := \langle \mathbf{s}, \mathbf{v}(f) \rangle = \frac{1}{m} \sum_{i=1}^m y_i \langle \mathbf{a}_i, \mathbf{v}(f) \rangle$.

To proceed, we use the following symmetrization bound, which is the complex-valued version of [35, Lemma 5.1].

Lemma 4. *Let $\{\epsilon_i\}_{i=1}^m$ be a sequence of independent complex-valued random variables, where $\epsilon_i \sim \epsilon = e^{j2\pi\theta}$, where θ uniformly distributed between $[0, 1)$. Then*

$$\begin{aligned} \mu &:= \mathbb{E} \left[\sup_{f \in [0,1]} |g_{\mathbf{x}^*}(f) - \mathbb{E}[g_{\mathbf{x}^*}(f)]| \right] \\ &\leq 2\mathbb{E} \left[\sup_{f \in [0,1]} \frac{1}{m} \left| \sum_{i=1}^m \epsilon_i y_i \langle \mathbf{a}_i, \mathbf{v}(f) \rangle \right| \right]. \end{aligned} \quad (\text{A.9})$$

Furthermore, we have the deviation inequality

$$\begin{aligned} &\mathbb{P} \left\{ \sup_{f \in [0,1]} |g_{\mathbf{x}^*}(f) - \mathbb{E}[g_{\mathbf{x}^*}(f)]| \geq 2\mu + t \right\} \\ &\leq 4\mathbb{P} \left\{ \sup_{f \in [0,1]} \frac{1}{m} \left| \sum_{i=1}^m \epsilon_i y_i \langle \mathbf{a}_i, \mathbf{v}(f) \rangle \right| > \frac{t}{2} \right\}. \end{aligned} \quad (\text{A.10})$$

Before applying Lemma 4, note that by symmetrization and rotational invariance, $\epsilon_i y_i \mathbf{a}_i$ have the same i.i.d. distribution of $\sqrt{2} \mathbf{a}_i$. Therefore, the following quantities are equivalent in distribution:

$$\begin{aligned} \sup_{f \in [0,1]} \frac{1}{m} \left| \sum_{i=1}^m \epsilon_i y_i \langle \mathbf{a}_i, \mathbf{v}(f) \rangle \right| &\sim \frac{\sqrt{2}}{m} \sup_{f \in [0,1]} \left| \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v}(f) \rangle \right| \\ &\sim \sqrt{\frac{2}{m}} \sup_{f \in [0,1]} |\langle \mathbf{g}, \mathbf{v}(f) \rangle|, \end{aligned}$$

where \mathbf{g} is a vector composed of i.i.d. $\mathcal{CN}(0, 1)$.

Applying (A.10) in Lemma 4 to (A.8), we have

$$\mathbb{P}(\|\mathbf{w}\|_{\mathcal{A}}^* \geq 2\mu + t) \leq 4\mathbb{P}\left(\sqrt{\frac{2}{m}} \sup_{f \in [0,1]} |\langle \mathbf{g}, \mathbf{v}(f) \rangle| \geq \frac{t}{2}\right). \quad (\text{A.11})$$

From (A.9) in Lemma 4, we have

$$\begin{aligned} \mu = \mathbb{E}[\|\mathbf{w}\|_{\mathcal{A}}^*] &\leq 2\sqrt{\frac{2}{m}} \mathbb{E}\left[\sup_{f \in [0,1]} |\langle \mathbf{g}, \mathbf{v}(f) \rangle|\right] \\ &\leq C\sqrt{\frac{n \log n}{m}}, \end{aligned} \quad (\text{A.12})$$

where the second line follows from [29, Appendix C,D] as

$$\mathbb{E}\left[\sup_{f \in [0,1]} |\langle \mathbf{g}, \mathbf{v}(f) \rangle|\right] \leq C_1 \sqrt{n \log(n)}.$$

Moreover, from [29, Appendix C], we have

$$\sup_{f \in [0,1]} |\langle \mathbf{g}, \mathbf{v}(f) \rangle| \leq C_2 \cdot \sqrt{n \log n}$$

hold with probability at least $1 - 1/(\pi n \log n)$. Set $t = 2C_2 \sqrt{n \log n}$ and plug in the above two inequalities in (A.11), we have that

$$\|\mathbf{w}\|_{\mathcal{A}}^* \leq C \cdot \sqrt{\frac{n \log n}{m}}$$

holds with probability at least $1 - 1/(\pi n \log n)$.

Appendix B: Proofs for Chapter 3

B.1 Gradient and Hessian of the Population Loss

For the convenience of analysis, we first provide the gradient and the Hessian formula for the cross-entropy loss using FCN and CNN here.

B.1.1 The FCN case

Consider the population loss function $f(\mathbf{W}) = \mathbb{E}[f_n(\mathbf{W})] = \mathbb{E}[\ell(\mathbf{W}; \mathbf{x})]$, where $\ell(\mathbf{W}; \mathbf{x})$ is associated with network $H_{\text{FCN}}(\mathbf{W}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}_k^\top \mathbf{x})$. Hiding the dependence on \mathbf{x} for notational simplicity, we can calculate the gradient and the Hessian as

$$\mathbb{E} \left[\frac{\partial \ell(\mathbf{W})}{\partial \mathbf{w}_j} \right] = \mathbb{E} \left[-\frac{1}{K} \frac{(y - H(\mathbf{W}))}{H(\mathbf{W})(1 - H(\mathbf{W}))} \phi'(\mathbf{w}_j^\top \mathbf{x}) \mathbf{x} \right], \quad (\text{B.1})$$

$$\mathbb{E} \left[\frac{\nabla^2 \ell(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_l} \right] = \mathbb{E} [\xi_{j,l}(\mathbf{W}) \cdot \mathbf{x} \mathbf{x}^\top], \quad (\text{B.2})$$

for $1 \leq j, l \leq K$. Here, when $j \neq l$,

$$\xi_{j,l}(\mathbf{W}) = \frac{1}{K^2} \phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x}) \cdot \frac{H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W})}{H^2(\mathbf{W})(1 - H(\mathbf{W}))^2},$$

and when $j = l$,

$$\begin{aligned}\xi_{j,j}(\mathbf{W}) &= \frac{1}{K^2} \phi'(\mathbf{w}_j^\top \mathbf{x})^2 \cdot \frac{H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W})}{H^2(\mathbf{W})(1 - H(\mathbf{W}))^2} \\ &\quad - \frac{1}{K} \phi''(\mathbf{w}_j^\top \mathbf{x}) \cdot \frac{y - H(\mathbf{W})}{H(\mathbf{W})(1 - H(\mathbf{W}))}.\end{aligned}$$

B.1.2 The CNN case

For the CNN case, i.e., $H(\mathbf{w}) := H_{\text{CNN}}(\mathbf{w}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}^\top \mathbf{x}^{(k)})$, the corresponding gradient and Hessian of the population loss function $\ell(\mathbf{w})$ is given by

$$\mathbb{E} \left[\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}} \right] = \mathbb{E} \left[-\phi'(\mathbf{w}^\top \mathbf{x}^{(1)}) \cdot \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \cdot \mathbf{x}^{(1)} \right], \quad (\text{B.3})$$

$$\mathbb{E} \left[\frac{\nabla^2 \ell(\mathbf{w})}{\partial \mathbf{w}^2} \right] = \mathbb{E} \left[\sum_{j=1}^K \sum_{l=1}^K g_{j,l}(\mathbf{w}) \mathbf{x}^{(j)} \mathbf{x}^{(l)\top} \right], \quad (\text{B.4})$$

where when $j \neq l$,

$$g_{j,l}(\mathbf{w}) = \frac{1}{K^2} \cdot \frac{H(\mathbf{w})^2 + y - 2y \cdot H(\mathbf{w})}{(H(\mathbf{w})(1 - H(\mathbf{w})))^2} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)}),$$

and when $j = l$,

$$\begin{aligned}g_{j,j}(\mathbf{w}) &= \frac{1}{K^2} \cdot \frac{H(\mathbf{w})^2 + y - 2y \cdot H(\mathbf{w})}{(H(\mathbf{w})(1 - H(\mathbf{w})))^2} \cdot \phi'(\mathbf{w}^\top \mathbf{x}^{(j)})^2 \\ &\quad - \frac{1}{K} \cdot \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \cdot \phi''(\mathbf{w}^\top \mathbf{x}^{(j)}).\end{aligned}$$

B.2 Proof of Theorem 3

In order to show that the empirical loss possesses a local strong convexity, we follow the following steps:

1. We first show that the Hessian $\nabla^2 f(\mathbf{W})$ of the population loss function is smooth with respect to $\nabla^2 f(\mathbf{W}^*)$ (Lemma 5);

2. We then show that $\nabla^2 f(\mathbf{W})$ satisfies local strong convexity and smoothness in a neighborhood of \mathbf{W}^* with appropriately chosen radius, $\mathbb{B}(\mathbf{W}^*, r)$, by leveraging similar properties of $\nabla^2 f(\mathbf{W}^*)$ (Lemma 6);
3. Next, we show that the Hessian of the empirical loss function $\nabla^2 f_n(\mathbf{W})$ is close to its population counterpart $\nabla^2 f(\mathbf{W})$ uniformly in $\mathbb{B}(\mathbf{W}^*, r)$ with high probability (Lemma 7).
4. Finally, putting all the arguments together, we establish $\nabla^2 f_n(\mathbf{W})$ satisfies local strong convexity and smoothness in $\mathbb{B}(\mathbf{W}^*, r)$.

To begin, we first show that the Hessian of the population risk is smooth enough around \mathbf{W}^* in the following lemmas.

Lemma 5 (Hessian Smoothness of Population Loss). *Suppose the loss $\ell(\cdot)$ associates with FCN (3.1), and assume $\|\mathbf{w}_k^*\|_2 \leq 1$ for all k and $\|\mathbf{W} - \mathbf{W}^*\|_{\text{F}} \leq 0.7$. Then we have*

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \leq \frac{C_1}{K^{\frac{3}{2}}} \cdot \|\mathbf{W} - \mathbf{W}^*\|_{\text{F}}, \quad (\text{B.5})$$

holds. Similarly, suppose the loss $\ell(\cdot)$ associates with CNN (3.2), and assume $\|\mathbf{w}^\|_2 \leq 1$ and $\|\mathbf{w} - \mathbf{w}^*\|_2 \leq 0.7$. We have*

$$\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| \leq C_2 \cdot K \cdot \|\mathbf{w} - \mathbf{w}^*\|_2, \quad (\text{B.6})$$

holds. Here C_1 and C_2 denote some large constants.

The proof is provided in Appendix B.4.1. Together with the fact that $\nabla^2 f(\mathbf{W}^*)$ be lower and upper bounded, Lemma 5 allows us to bound $\nabla^2 f(\mathbf{W})$ in a neighborhood around ground truth, given below.

Lemma 6 (Local Strong Convexity and Smoothness of Population Loss). *If the loss $\ell(\cdot)$ associates with FCN (3.1), there exists some constant C_1 , such that*

$$\frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{W}) \preceq C_1 \cdot \mathbf{I},$$

holds for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^, r_{\text{FCN}})$ with $r_{\text{FCN}} := \frac{C_2}{K^{\frac{1}{2}}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}$. Moreover, if loss $\ell(\cdot)$ associates with CNN (3.2), then we have*

$$C_3 \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{w}) \preceq C_4 \cdot K \cdot \mathbf{I}, \quad (\text{B.7})$$

holds for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^, r_{\text{CNN}})$ with $r_{\text{CNN}} := C_5 \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K^2}$.*

The proof is provided in Appendix B.4.2. The next step is to show the Hessian of the empirical loss function is close to the Hessian of the population loss function in a uniform sense, which can be summarized as follows.

Lemma 7. *For the loss $\ell(\cdot)$ associated with FCN (3.1), there exists a constant C such that as long as $n \geq C \cdot dK \log dK$, with probability at least $1 - d^{-10}$, the following holds*

$$\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})} \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| \leq C \sqrt{\frac{dK \log n}{n}}, \quad (\text{B.8})$$

where $r_{\text{FCN}} := \frac{C}{K^{\frac{1}{2}}} \cdot \frac{\rho(\sigma_K)}{\kappa^2 \lambda}$. For the loss $\ell(\cdot)$ associated with CNN (3.2), we have

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})} \|\nabla^2 f_n(\mathbf{w}) - \nabla^2 f(\mathbf{w})\| \leq CK^2 \sqrt{\frac{\frac{d}{K} \cdot \log(n)}{n}}, \quad (\text{B.9})$$

holds with probability at least $1 - d^{-10}$, as long as $n \geq \frac{d}{K} \log\left(\frac{d}{K}\right)$, and $r_{\text{CNN}} := C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^\|_2)}{K^2}$.*

The proof is provided in Appendix B.4.3. Combining the above results will give us the result. Next we assume that the loss $\ell(\cdot)$ associates with FCN, and take it as

an example in the proof. Then if the loss $\ell(\cdot)$ associates with CNN, the proof follows in the same manner.

Proof of Theorem 3. With probability at least $1 - d^{-10}$,

$$\begin{aligned}\nabla^2 f_n(\mathbf{W}) &\succeq \nabla^2 f(\mathbf{W}) - \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| \cdot \mathbf{I} \\ &\succeq \Omega\left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}\right) \cdot \mathbf{I} - \Omega\left(C \cdot \sqrt{\frac{dK \log n}{n}}\right) \cdot \mathbf{I}.\end{aligned}$$

As long as the sample size n is set to satisfy

$$C \cdot \sqrt{\frac{dK \log n}{n}} \leq \frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda},$$

i.e. $n \geq C \cdot dK^5 \log^2 d \cdot \left(\frac{\kappa^2 \lambda}{\rho_{\text{FCN}}(\sigma_K)}\right)^2$, we have

$$\nabla^2 f_n(\mathbf{W}) \succeq \Omega\left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}\right) \cdot \mathbf{I}.$$

holds for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$. Similarly, we have

$$\nabla^2 f_n(\mathbf{W}) \preceq C \cdot \mathbf{I}$$

holds for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$. □

B.3 Proof of Theorem 4

We have established that $f_n(\mathbf{W})$ is strongly convex in $\mathbb{B}(\mathbf{W}^*, r)$ in Theorem 3. Thus there exists at most one critical point in $\mathbb{B}(\mathbf{W}^*, r)$. The proof of Theorem 4 follows the steps below:

1. We first show that the gradient $\nabla f_n(\mathbf{W})$ concentrates around $\nabla f(\mathbf{W})$ in $\mathbb{B}(\mathbf{W}^*, r)$ (Lemma 8), and then invoke [50, Theorem 2] to guarantee that there indeed exists a critical point $\widehat{\mathbf{W}}_n$ in $\mathbb{B}(\mathbf{W}^*, r)$;

2. We next show that $\widehat{\mathbf{W}}_n$ is close to \mathbf{W}^* and gradient descent converges linearly to $\widehat{\mathbf{W}}_n$ with a properly chosen step size.

To begin, the following lemma establishes that $\nabla f_n(\mathbf{W})$ uniformly concentrates around $\nabla f(\mathbf{W})$.

Lemma 8. *If the loss $\ell(\cdot)$ associates with FCN (3.1) with $r_{\text{FCN}} := \frac{C}{K^{\frac{1}{2}}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}$, and $\|\mathbf{w}_k^*\|_2 \leq 1$ for all k , then*

$$\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})} \|\nabla f_n(\mathbf{W}) - \nabla f(\mathbf{W})\| \leq C \sqrt{\frac{d\sqrt{K} \log n}{n}}$$

holds with probability at least $1 - d^{-10}$, as long as $n \geq CdK \log(dK)$. If the loss $\ell(\cdot)$ associates with CNN (3.2), with $r_{\text{CNN}} := C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K^2}$ and $\|\mathbf{w}^*\|_2 \leq 1$, then

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})} \|\nabla f_n(\mathbf{w}) - \nabla f(\mathbf{w})\| \leq C \cdot \sqrt{\frac{d \log n}{n}} \quad (\text{B.10})$$

holds with probability at least $1 - d^{-10}$ as long as $n \geq C \frac{d}{K} \log\left(\frac{d}{K}\right)$.

The proof is provided in Appendix B.4.4. Notice that for the population risk function $f(\mathbf{W})$, \mathbf{W}^* is the unique critical point in $\mathbb{B}(\mathbf{W}^*, r)$ due to local strong convexity. With Lemma 7 and Lemma 8, we can invoke [50, Theorem 2], which guarantees the following.

Corollary 1. *If the loss $\ell(\cdot)$ associates with FCN or CNN, there exists one and only one critical point $\widehat{\mathbf{W}}_n \in \mathbb{B}(\mathbf{W}^*, r)$ that satisfies $\nabla f_n(\widehat{\mathbf{W}}_n) = \mathbf{0}$ correspondingly.*

Again, since the proof for the case with the loss $\ell(\cdot)$ associating with FCN is the same as that for CNN, we next take FCN as an example.

We first show that $\widehat{\mathbf{W}}_n$ is close to \mathbf{W}^* . By the extended mean value theorem, there exists \mathbf{W}' on the straight line connecting \mathbf{W}^* and $\widehat{\mathbf{W}}_n$ such that

$$\begin{aligned} f_n(\widehat{\mathbf{W}}_n) &= f_n(\mathbf{W}^*) + \left\langle \nabla f_n(\mathbf{W}^*), \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*) \right\rangle \\ &\quad + \frac{1}{2} \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*)^\top \nabla^2 f_n(\mathbf{W}') \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*) \\ &\leq f_n(\mathbf{W}^*), \end{aligned} \tag{B.11}$$

where the last inequality follows from the optimality of $\widehat{\mathbf{W}}_n$. By Theorem 3, we have

$$\begin{aligned} &\frac{1}{2} \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*)^\top \nabla^2 f_n(\mathbf{W}') \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*) \\ &\geq \Omega \left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \right) \left\| \widehat{\mathbf{W}}_n - \mathbf{W}^* \right\|_{\text{F}}^2. \end{aligned} \tag{B.12}$$

On the other hand, by the Cauchy-Schwarz inequality, we have

$$\begin{aligned} \left| \left\langle \nabla f_n(\mathbf{W}^*), \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*) \right\rangle \right| &\leq \|\nabla f_n(\mathbf{W}^*)\|_2 \|\widehat{\mathbf{W}}_n - \mathbf{W}^*\|_{\text{F}} \\ &\leq \Omega \left(\sqrt{\frac{dK^{1/2} \log n}{n}} \right) \|\widehat{\mathbf{W}}_n - \mathbf{W}^*\|_{\text{F}}, \end{aligned} \tag{B.13}$$

where the last line follows from Lemma 8. Plugging (B.12) and (B.13) into (B.11), we have

$$\|\widehat{\mathbf{W}}_n - \mathbf{W}^*\|_{\text{F}} \leq \Omega \left(\frac{K^{\frac{9}{4}} \kappa^2 \lambda}{\rho_{\text{FCN}}(\sigma_K)} \sqrt{\frac{d \log n}{n}} \right). \tag{B.14}$$

Now we have established that there indeed exists a critical point in $\mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$.

We can then establish the local linear convergence of gradient descent as below. Let \mathbf{W}_t be the estimate at the t -th iteration. Due to the update rule, we have

$$\begin{aligned} \mathbf{W}_{t+1} - \widehat{\mathbf{W}}_n &= \mathbf{W}_t - \eta \nabla f_n(\mathbf{W}_t) - \left(\widehat{\mathbf{W}}_n - \eta \nabla f_n(\widehat{\mathbf{W}}_n) \right) \\ &= \left(\mathbf{I} - \eta \int_0^1 \nabla^2 f_n(\mathbf{W}(\gamma)) \right) \left(\mathbf{W}_t - \widehat{\mathbf{W}}_n \right), \end{aligned}$$

where $\mathbf{W}(\gamma) = \widehat{\mathbf{W}}_n + \gamma (\mathbf{W}_t - \widehat{\mathbf{W}}_n)$ for $\gamma \in [0, 1]$. If $\mathbf{W}_t \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$, it is obvious that $\mathbf{W}(\gamma) \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$, and by Theorem 3, we have

$$H_{\min} \cdot \mathbf{I} \preceq \nabla^2 f_n(\mathbf{W}(\gamma)) \preceq H_{\max} \cdot \mathbf{I},$$

where $H_{\min} = \Omega\left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}\right)$ and $H_{\max} = C$. Therefore, we have

$$\begin{aligned} \|\mathbf{W}_{t+1} - \widehat{\mathbf{W}}_n\|_{\text{F}} &\leq \|\mathbf{I} - \eta \int_0^1 \nabla^2 f_n(\mathbf{W}(\gamma))\| \|\mathbf{W}_t - \widehat{\mathbf{W}}_n\|_{\text{F}} \\ &\leq (1 - \eta H_{\min}) \|\mathbf{W}_t - \widehat{\mathbf{W}}_n\|_{\text{F}}. \end{aligned} \quad (\text{B.15})$$

Hence, by setting $\eta = \frac{1}{H_{\max}} := \Omega(C)$, we obtain

$$\|\mathbf{W}_{t+1} - \widehat{\mathbf{W}}_n\|_{\text{F}} \leq \left(1 - \frac{H_{\min}}{H_{\max}}\right) \|\mathbf{W}_t - \widehat{\mathbf{W}}_n\|_{\text{F}}, \quad (\text{B.16})$$

which implies that gradient descent converges linearly to the local minimizer $\widehat{\mathbf{W}}_n$.

B.4 Proof of Auxiliary Lemmas

B.4.1 Proof of Lemma 5.

We prove the two claims for FCN and CNN separately as below.

- **The FCN case:** Let $\Delta = \nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)$. For each $(j, l) \in [K] \times [K]$, let $\Delta_{j,l} \in \mathbb{R}^{d \times d}$ denote the (j, l) -th block of Δ . Let $\mathbf{a} = [\mathbf{a}_1^\top, \dots, \mathbf{a}_K^\top]^\top \in \mathbb{R}^{dK}$.

By definition,

$$\begin{aligned} \|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| &= \max_{\|\mathbf{a}\|=1} \mathbf{a}^\top (\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)) \mathbf{a} \\ &= \max_{\|\mathbf{a}\|=1} \sum_{j=1}^K \sum_{l=1}^K \mathbf{a}_j^\top \Delta_{j,l} \mathbf{a}_l. \end{aligned} \quad (\text{B.17})$$

From (B.2) we know that

$$\Delta_{j,l} = \frac{\partial^2 f(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_l} - \frac{\partial^2 f(\mathbf{W}^*)}{\partial \mathbf{w}_j^* \partial \mathbf{w}_l^*} = \mathbb{E} [(\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}^*)) \cdot \mathbf{x} \mathbf{x}^\top], \quad (\text{B.18})$$

and then by the mean value theorem, we can further expand $\xi_{j,l}(\mathbf{W})$ as

$$\xi_{j,l}(\mathbf{W}) = \xi_{j,l}(\mathbf{W}^*) + \sum_{k=1}^K \left\langle \frac{\partial \xi_{j,l}(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{w}}_k}, \mathbf{w}_k - \mathbf{w}_k^* \right\rangle, \quad (\text{B.19})$$

where $\tilde{\mathbf{W}} = \eta \cdot \mathbf{W} + (1 - \eta) \mathbf{W}^*$ for some $\eta \in (0, 1)$. Thus we can write $\Delta_{j,l}$ as

$$\Delta_{j,l} = \mathbb{E} \left[\left(\sum_{k=1}^K \left\langle \frac{\partial \xi_{j,l}(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{w}}_k}, \mathbf{w}_k - \mathbf{w}_k^* \right\rangle \right) \cdot \mathbf{x} \mathbf{x}^\top \right], \quad (\text{B.20})$$

which can be further simplified as

$$\Delta_{j,l} = \mathbb{E} \left[\left(\sum_{k=1}^K T_{j,l,k} \langle \mathbf{x}, \mathbf{w}_k - \mathbf{w}_k^* \rangle \right) \cdot \mathbf{x} \mathbf{x}^\top \right], \quad (\text{B.21})$$

by the fact that $\frac{\partial \xi_{j,l}(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{w}}_k}$ can be written as $T_{j,l,k} \cdot \mathbf{x}$, where $T_{j,l,k} \in \mathbb{R}$ is a scalar depending on \mathbf{x} . When $j = l$, we calculate $\frac{\partial \xi_{j,l}(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{w}}_k}$ for illustration,

$$\begin{aligned} \frac{\partial \xi_{j,j}(\mathbf{W})}{\partial \mathbf{w}_k} &= \left(-\frac{2}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})^2}{H(\mathbf{W})^3} + \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})^2} \right) \frac{1}{K} \phi(\mathbf{w}_k^\top \mathbf{x}) \mathbf{x}, \quad k \neq j \\ \frac{\partial \xi_{j,j}(\mathbf{W})}{\partial \mathbf{w}_k} &= \left(\frac{2}{K^2} \left(\frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})^2} - \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})^2}{H(\mathbf{W})^3} \right) \right. \\ &\quad \left. + \frac{1}{K} \left(\frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})^2} - \frac{\phi'''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})} \right) \right) \frac{1}{K} \phi(\mathbf{w}_k^\top \mathbf{x}) \mathbf{x}, \quad k = j \end{aligned} \quad (\text{B.22})$$

where we have simplified the presentation by setting $y = 1$, since y is a binary random variable, and we will show that in either case $|T_{j,j,k}|$ is upper bounded, i.e., in this case

$$|T_{j,j,k}| \leq \begin{cases} \max \left\{ \frac{2}{K^3} \frac{1}{H(\tilde{\mathbf{W}})^3}, \frac{1}{K^2} \frac{1}{H(\tilde{\mathbf{W}})^2} \right\} & y = 1 \\ \max \left\{ \frac{2}{K^3} \frac{1}{(1-H(\tilde{\mathbf{W}}))^3}, \frac{1}{K^2} \frac{1}{(1-H(\tilde{\mathbf{W}}))^2} \right\} & y = 0 \end{cases},$$

since $\phi(\cdot), \phi'(\cdot), \phi''(\cdot), \phi'''(\cdot)$ are bounded. More generally, by calculating the other case we can claim that

$$|T_{j,l,k}| \leq \max \left\{ \frac{2}{K^3} \frac{1}{H(\tilde{\mathbf{W}})^3}, \frac{1}{K^2} \frac{1}{H(\tilde{\mathbf{W}})^2}, \frac{2}{K^3} \frac{1}{(1-H(\tilde{\mathbf{W}}))^3}, \frac{1}{K^2} \frac{1}{(1-H(\tilde{\mathbf{W}}))^2} \right\}, \quad (\text{B.23})$$

holds for all j, l, k . Then, we can upper bound $\mathbf{a}_j^\top \Delta_{j,l} \mathbf{a}_l$ using Cauchy-Schwarz inequality,

$$\begin{aligned}
\mathbf{a}_j^\top \Delta_{j,l} \mathbf{a}_l &= \mathbb{E} \left[\left(\sum_{k=1}^K T_{j,l,k} \langle \mathbf{x}, \mathbf{w}_k - \mathbf{w}_k^* \rangle \right) \cdot (\mathbf{a}_j^\top \mathbf{x}) (\mathbf{a}_l^\top \mathbf{x}) \right] \\
&\leq \sqrt{\mathbb{E} \left[\sum_{k=1}^K T_{j,l,k}^2 \right]} \cdot \sqrt{\mathbb{E} \left[\sum_{k=1}^K (\langle \mathbf{x}, \mathbf{w}_k - \mathbf{w}_k^* \rangle (\mathbf{a}_j^\top \mathbf{x}) (\mathbf{a}_l^\top \mathbf{x}))^2 \right]} \\
&\leq \sqrt{\sum_{k=1}^K \mathbb{E} [T_{j,l,k}^2]} \cdot \sqrt{\sum_{k=1}^K \|\mathbf{w}_k - \mathbf{w}_k^*\|_2^2 \cdot \|\mathbf{a}_j\|_2^2 \cdot \|\mathbf{a}_l\|_2^2}. \quad (\text{B.24})
\end{aligned}$$

Plug it back to (B.17) we can obtain the following inequality,

$$\begin{aligned}
&\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \\
&\leq \max_{\|\mathbf{a}\|=1} \sum_{j=1}^K \sum_{l=1}^K \sqrt{\sum_{k=1}^K \mathbb{E} [T_{j,l,k}^2]} \cdot \sqrt{\sum_{k=1}^K \|\mathbf{w}_k - \mathbf{w}_k^*\|_2^2 \cdot \|\mathbf{a}_j\|_2^2 \cdot \|\mathbf{a}_l\|_2^2}. \quad (\text{B.25})
\end{aligned}$$

Then the problem boils down to upper bound $\mathbb{E} [T_{i,j,k}^2]$, which we can apply the following lemma, whose proof can be found in Section B.4.5.

Lemma 9. *Let $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $t = \max \{\|\mathbf{w}_1\|_2, \dots, \|\mathbf{w}_K\|_2\}$ and $z \in \mathbb{Z}$ such that $z \geq 1$, for the sigmoid activation function $\phi(x) = \frac{1}{1+e^{-x}}$, the following*

$$\begin{aligned}
\mathbb{E} \left[\left(\frac{1}{\frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_j^\top \mathbf{x})} \right)^z \right] &\leq C_1 \cdot e^{t^2}, \\
\mathbb{E} \left[\left(\frac{1}{\left(1 - \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_j^\top \mathbf{x})\right)} \right)^z \right] &\leq C_2 \cdot e^{t^2} \quad (\text{B.26})
\end{aligned}$$

holds for some large enough constants C_1, C_2 that depend on the constant z .

Setting $z = 4$ and $z = 6$ in Lemma 9, together with (B.23) we obtain that

$$\mathbb{E} [T_{j,l,k}^2] \leq \frac{C}{K^4} \cdot e^{\max_{1 \leq i \leq k} \|\tilde{\mathbf{w}}_i\|_2^2}, \quad (\text{B.27})$$

holds for some constant C . Plugging (B.27) into (B.25), we obtain

$$\begin{aligned} \|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| &\leq \frac{C}{K^{\frac{3}{2}}} e^{\|\tilde{\mathbf{w}}\|_F^2} \cdot \|\mathbf{W} - \mathbf{W}^*\|_F \cdot \max_{\|\mathbf{a}\|=1} \sum_{j=1}^K \sum_{l=1}^K \|\mathbf{a}_j\|_2 \|\mathbf{a}_l\|_2 \\ &\leq \frac{C}{K^{\frac{3}{2}}} e^{\|\tilde{\mathbf{w}}\|_F^2} \cdot \|\mathbf{W} - \mathbf{W}^*\|_F. \end{aligned} \quad (\text{B.28})$$

Further since $e^{\max_{1 \leq i \leq k} \|\tilde{\mathbf{w}}_i\|_2^2} \leq C$ gives that $\|\mathbf{w}_i - \mathbf{w}_i^*\|_2 \leq 0.7$, where we have used the assumption that $\max_{1 \leq i \leq k} \|\mathbf{w}_i^*\|_2^2 \leq 1$, we conclude that

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \leq \frac{C}{K^{\frac{3}{2}}} \|\mathbf{W} - \mathbf{W}^*\|_F \quad (\text{B.29})$$

holds for some constant C .

- **The CNN case:** according to (B.4), we can calculate the upper bound of $\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\|$ by definition as

$$\max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \mathbb{E} [(g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}^*)) \cdot \mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u}]. \quad (\text{B.30})$$

We then again apply the mean value theorem to $g_{j,l}(\mathbf{w})$, such that there exists $\tilde{\mathbf{w}} = \eta \mathbf{w} + (1 - \eta) \mathbf{w}^*$ for some $\eta \in (0, 1)$,

$$g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}^*) = \langle \nabla g_{j,l}(\tilde{\mathbf{w}}), \mathbf{w} - \mathbf{w}^* \rangle.$$

Similarly to the FCN case, we can write $\nabla g_{j,l}(\tilde{\mathbf{w}})$ in the form of

$$\nabla g_{j,l}(\tilde{\mathbf{w}}) = \sum_{k=1}^K S_{j,l,k} \cdot \mathbf{x}^{(k)},$$

where $S_{j,l,k}$ is a scalar that depends on $\tilde{\mathbf{w}}$ and $\mathbf{x}^{(k)}$, $k = 1, \dots, K$. Again we take $j \neq l$ as an example to calculate $S_{j,l,k}$, by definition, and obtain

$$\begin{aligned} K^2 \cdot \frac{\partial g_{j,l}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{(1 - H(\mathbf{w})) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi''(\mathbf{w}^\top \mathbf{x}^{(l)})}{(1 - H(\mathbf{w}))^3} \cdot \mathbf{x}^{(l)} \\ &\quad + \frac{(1 - H(\mathbf{w})) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \phi''(\mathbf{w}^\top \mathbf{x}^{(j)})}{(1 - H(\mathbf{w}))^3} \cdot \mathbf{x}^{(j)} \\ &\quad - \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{(1 - H(\mathbf{w}))^3} \cdot \left(\frac{1}{K} \sum_{k=1}^K \mathbf{x}^{(k)} \right), \end{aligned} \quad (\text{B.31})$$

where we set $y = 0$ for simplification. Then we obtain

$$S_{j,l,l} = \frac{1}{K^2} \frac{(1 - H(\mathbf{w})) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi''(\mathbf{w}^\top \mathbf{x}^{(l)})}{(1 - H(\mathbf{w}))^3} - \frac{1}{K^3} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{(1 - H(\mathbf{w}))^3}. \quad (\text{B.32})$$

and

$$|S_{j,l,l}| \leq \frac{1}{K^2} \frac{1}{(1 - H(\tilde{\mathbf{w}}))^3}, \quad (\text{B.33})$$

hold, where we used the fact that $0 \leq H(\mathbf{w}) \leq 1$ and $\phi'(\cdot), \phi''(\cdot)$ are bounded.

Hence in the same way, we can obtain

$$|S_{j,l,k}| \leq \begin{cases} \max \left\{ \frac{1}{K^2} \frac{1}{(1 - H(\tilde{\mathbf{w}}))^3}, \frac{1}{K^2} \frac{1}{(H(\tilde{\mathbf{w}}))^3} \right\} & j \neq l \\ \max \left\{ \frac{1}{K} \frac{1}{(1 - H(\tilde{\mathbf{w}}))^2}, \frac{1}{K} \frac{1}{(H(\tilde{\mathbf{w}}))^2} \right\} & j = l \end{cases}. \quad (\text{B.34})$$

Plug these back to (B.30) we obtain

$$\begin{aligned} & \|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| \\ & \leq \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \mathbb{E} \left[\sum_{k=1}^K \langle S_{j,l,k} \cdot \mathbf{x}^{(k)}, \mathbf{w} - \mathbf{w}^* \rangle \cdot \mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u} \right] \\ & = \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \mathbb{E} \left[\sum_{k=1}^K S_{j,l,k} \cdot (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{x}^{(k)} \cdot \mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u} \right] \\ & \leq \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \sqrt{\mathbb{E} \left[\sum_{k=1}^K S_{j,l,k}^2 \right] \cdot \mathbb{E} \left[\sum_{k=1}^K \left((\mathbf{w} - \mathbf{w}^*)^\top \mathbf{x}^{(k)} \right)^2 (\mathbf{u}^\top \mathbf{x}^{(j)})^2 (\mathbf{x}^{(l)\top} \mathbf{u})^2 \right]} \\ & \leq \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \sqrt{\mathbb{E} \left[\sum_{k=1}^K S_{j,l,k}^2 \right] \cdot \sum_{k=1}^K \|\mathbf{w} - \mathbf{w}^*\|_2^2 \cdot \|\mathbf{u}\|_2^2 \cdot \|\mathbf{u}\|_2^2} \\ & \leq C \cdot K \cdot e^{\|\tilde{\mathbf{w}}\|_2^2} \cdot \|\mathbf{w} - \mathbf{w}^*\|_2, \end{aligned} \quad (\text{B.35})$$

where the second inequality follows from Cauchy-Schwarz inequality, and the

last inequality follows from (B.34) and Lemma 9. Further since $e^{\|\tilde{\mathbf{w}}\|_2^2} \leq C \cdot$

$(1 + \|\mathbf{w} - \mathbf{w}^*\|_2^2)$ given that $\|\mathbf{w} - \mathbf{w}^*\|_2 \leq 0.7$, we conclude that

$$\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| \leq C \cdot K \cdot \|\mathbf{w} - \mathbf{w}^*\|_2 \quad (\text{B.36})$$

holds for some constant C and $\|\mathbf{w} - \mathbf{w}^*\| \leq 0.7$.

B.4.2 Proof of Lemma 6

We first present upper and lower bounds on the Hessian $\nabla^2 f(\mathbf{W}^*)$ of the population risk at ground truth, and then apply Lemma 5 to obtain a uniform bound in the neighborhood of \mathbf{W}^* .

- **The FCN case:** Recall

$$\begin{aligned}\frac{\partial^2 f(\mathbf{W}^*)}{\partial \mathbf{w}_j^2} &= \mathbb{E} \left[\frac{1}{K^2} \cdot \left(\frac{\phi'(\mathbf{w}_j^{*\top} \mathbf{x})^2}{H(\mathbf{W}^*)(1-H(\mathbf{W}^*))} \right) \mathbf{x} \mathbf{x}^\top \right], \\ \frac{\partial^2 f(\mathbf{W}^*)}{\partial \mathbf{w}_j \partial \mathbf{w}_l} &= \mathbb{E} \left[\frac{1}{K^2} \cdot \left(\frac{\phi'(\mathbf{w}_j^{*\top} \mathbf{x}) \phi'(\mathbf{w}_l^{*\top} \mathbf{x})}{H(\mathbf{W}^*)(1-H(\mathbf{W}^*))} \right) \mathbf{x} \mathbf{x}^\top \right],\end{aligned}$$

where we have applied the fact that $\mathbb{E}[y|\mathbf{x}] = H(\mathbf{W}^*)$. Let $\mathbf{a} = [\mathbf{a}_1^\top, \dots, \mathbf{a}_K^\top]^\top \in \mathbb{R}^{dK}$. Then we can write

$$\begin{aligned}\nabla^2 f(\mathbf{W}^*) &\succeq \left(\min_{\|\mathbf{a}\|_2=1} \mathbf{a}^\top \nabla^2 f(\mathbf{W}^*) \mathbf{a} \right) \cdot \mathbf{I} \\ &= \min_{\|\mathbf{a}\|_2=1} \frac{1}{K^2} \mathbb{E} \left[\frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x}) \right)^2}{H(\mathbf{W}^*)(1-H(\mathbf{W}^*))} \right] \cdot \mathbf{I}.\end{aligned}\tag{B.37}$$

Since $0 \leq H(\mathbf{W}^*) \leq 1$, we have that $H(\mathbf{W}^*)(1-H(\mathbf{W}^*)) \leq \frac{1}{4}$. Hence,

$$\nabla^2 f(\mathbf{W}^*) \succeq \min_{\|\mathbf{a}\|_2=1} \frac{4}{K^2} \mathbb{E} \left[\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x}) \right)^2 \right] \cdot \mathbf{I} \succeq \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \cdot \mathbf{I},\tag{B.38}$$

where the last inequality follows from [7, Lemmas D.4 and D.6]. To derive an upper bound of $\nabla^2 f(\mathbf{W}^*)$, we have

$$\begin{aligned}\nabla^2 f(\mathbf{W}^*) &\preceq \left(\max_{\|\mathbf{a}\|_2=1} \mathbf{a}^\top \nabla^2 f(\mathbf{W}^*) \mathbf{a} \right) \cdot \mathbf{I} \\ &= \max_{\|\mathbf{a}\|_2=1} \frac{1}{K^2} \mathbb{E} \left[\frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x}) \right)^2}{\frac{1}{K^2} \sum_{j,l} \phi(\mathbf{w}_j^{*\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{*\top} \mathbf{x}))} \right].\end{aligned}\tag{B.39}$$

Then by Cauchy-Schwarz inequality, we have

$$\frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{\star\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x})\right)^2}{\frac{1}{K^2} \sum_{j,l} \phi(\mathbf{w}_j^{\star\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{\star\top} \mathbf{x}))} \leq \frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{\star\top} \mathbf{x})\right)^2 \cdot \left(\sum_{j=1}^K (\mathbf{a}_j^\top \mathbf{x})^2\right)}{\frac{1}{K^2} \sum_{j,l} \phi(\mathbf{w}_j^{\star\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{\star\top} \mathbf{x}))}.$$

Further since $\phi'(\mathbf{w}_j^{\star\top} \mathbf{x}) \leq \frac{1}{4}$, and

$$\begin{aligned} \sum_{j,l} \phi(\mathbf{w}_j^{\star\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{\star\top} \mathbf{x})) &\geq \sum_{j=1}^K \phi(\mathbf{w}_j^{\star\top} \mathbf{x}) (1 - \phi(\mathbf{w}_j^{\star\top} \mathbf{x})) \\ &= \sum_{j=1}^K \phi'(\mathbf{w}_j^{\star\top} \mathbf{x}) \\ &\geq 4 \sum_{j=1}^K \phi'(\mathbf{w}_j^{\star\top} \mathbf{x})^2, \end{aligned} \quad (\text{B.40})$$

we obtain

$$\mathbf{a}^\top \nabla^2 f(\mathbf{W}^*) \mathbf{a} \leq \frac{1}{K^2} \mathbb{E} \left[\frac{CK^2}{4} \sum_{j=1}^K (\mathbf{a}_j^\top \mathbf{x})^2 \right]. \quad (\text{B.41})$$

Plugging (B.41) back to (B.39), we obtain

$$\nabla^2 f(\mathbf{W}^*) \preceq C \cdot \mathbf{I}. \quad (\text{B.42})$$

Thus together with the lower bound (B.38), we conclude that

$$\frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{W}^*) \preceq C \cdot \mathbf{I}. \quad (\text{B.43})$$

From Lemma 5, we have

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \lesssim \frac{C}{K^{\frac{3}{2}}} \|\mathbf{W} - \mathbf{W}^*\|_F. \quad (\text{B.44})$$

Therefore, if $\|\mathbf{W}^* - \mathbf{W}\|_F \leq 0.7$ and

$$\frac{C}{K^{\frac{3}{2}}} \cdot \|\mathbf{W} - \mathbf{W}^*\|_F \leq \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda},$$

i.e., if $\|\mathbf{W} - \mathbf{W}^*\|_F \leq \min \left\{ \frac{C}{K^{\frac{1}{2}}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}, 0.7 \right\}$ for some constant C , we have

$$\begin{aligned} \sigma_{\min}(\nabla^2 f(\mathbf{W})) &\geq \sigma_{\min}(\nabla^2 f(\mathbf{W}^*)) - \|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \\ &\gtrsim \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} - \frac{C}{K^{\frac{3}{2}}} \|\mathbf{W} - \mathbf{W}^*\|_F \\ &\gtrsim \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}. \end{aligned}$$

Moreover, within the same neighborhood, by the triangle inequality we have

$$\|\nabla^2 f(\mathbf{W})\| \leq \|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| + \|\nabla^2 f(\mathbf{W}^*)\| \lesssim C.$$

- **The CNN case:** Following from (B.4), we have

$$\nabla^2 f(\mathbf{w}^*) = \mathbb{E} \left[\frac{\frac{1}{K^2} \sum_{j,l} \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(j)}) \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(l)}) \mathbf{x}^{(j)} \mathbf{x}^{(l)\top}}{H(\mathbf{w}^*) (1 - H(\mathbf{w}^*))} \right]. \quad (\text{B.45})$$

By definition, we lower bound $\nabla^2 f(\mathbf{w}^*)$ by

$$\begin{aligned} &\min_{\|\mathbf{u}\|=1} \mathbb{E} \left[\frac{\frac{1}{K^2} \sum_{j,l} \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)} \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(l)}) \mathbf{u}^\top \mathbf{x}^{(l)}}{H(\mathbf{w}^*) (1 - H(\mathbf{w}^*))} \right] \cdot \mathbf{I} \\ &\succeq \min_{\|\mathbf{u}\|=1} \mathbb{E} \left[\frac{4}{K^2} \sum_{j,l} \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)} \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(l)}) \mathbf{u}^\top \mathbf{x}^{(l)} \right] \cdot \mathbf{I} \\ &= \frac{4}{K^2} \cdot \left(\min_{\|\mathbf{u}\|=1} \sum_{j \neq l} \mathbb{E} [\phi'(\mathbf{w}^{*\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)}] \cdot \mathbb{E} [\phi'(\mathbf{w}^{*\top} \mathbf{x}^{(l)}) \mathbf{u}^\top \mathbf{x}^{(l)}] \right. \\ &\quad \left. + \sum_{j=1}^K \mathbb{E} [(\phi'(\mathbf{w}^{*\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)})^2] \right) \cdot \mathbf{I}, \end{aligned}$$

where the last equality follows from the fact that $\mathbf{x}^{(j)}$ is independent from $\mathbf{x}^{(l)}$ given that $j \neq l$. Next we decompose \mathbf{u} as $\mathbf{u} = \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2} \cdot \mathbf{w}^* + \left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2} \cdot \mathbf{w}^* \right)$,

and calculate the expectation as

$$\begin{aligned}
& \mathbb{E} [\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)}] \\
&= \mathbb{E} \left[\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \left(\frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \cdot \mathbf{w}^{\star} + \left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \cdot \mathbf{w}^{\star} \right) \right)^\top \mathbf{x}^{(j)} \right] \\
&= \mathbb{E} \left[\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \cdot \mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right] \\
&\quad + \mathbb{E} [\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)})] \cdot \mathbb{E} \left[\left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \cdot \mathbf{w}^{\star} \right)^\top \mathbf{x}^{(j)} \right] \\
&= \frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \mathbb{E} [\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{w}^{\star\top} \mathbf{x}^{(j)}],
\end{aligned}$$

where the second equality follows from the independence of $\mathbf{w}^{\star\top} \mathbf{x}^{(j)}$ and

$\left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \cdot \mathbf{w}^{\star} \right)^\top \mathbf{x}^{(j)}$. Hence,

$$\begin{aligned}
& \mathbb{E} [\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)}] \cdot \mathbb{E} [\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(l)}) \mathbf{u}^\top \mathbf{x}^{(l)}] \\
&= \left(\frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \right)^2 (\mathbb{E} [\phi' (z) z])^2 \\
&= 0,
\end{aligned}$$

where $z = \mathbf{w}^{\star\top} \mathbf{x}^{(j)} \sim \mathcal{N}(0, \|\mathbf{w}^{\star}\|_2^2)$, and the last equality follows because

$\phi' (z) z = -(\phi' (-z) \cdot (-z))$. Similarly,

$$\begin{aligned}
& \mathbb{E} \left[(\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)})^2 \right] \\
&= \mathbb{E} \left[\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)})^2 \cdot \left(\left(\frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \cdot \mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right)^2 + \left(\left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \cdot \mathbf{w}^{\star} \right)^\top \mathbf{x}^{(j)} \right)^2 \right) \right] \\
&= \left(\frac{\mathbf{u}^\top \mathbf{w}^{\star}}{\|\mathbf{w}^{\star}\|_2^2} \right)^2 \cdot \mathbb{E} \left[\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)})^2 (\mathbf{w}^{\star\top} \mathbf{x}^{(j)})^2 \right] \\
&\quad + \left(\|\mathbf{u}\|_2^2 - \frac{(\mathbf{u}^\top \mathbf{w}^{\star})^2}{\|\mathbf{w}^{\star}\|_2^2} \right) \cdot \mathbb{E} \left[\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)})^2 \right]. \tag{B.46}
\end{aligned}$$

Together with Definition 2, we have

$$\mathbb{E} \left[(\phi' (\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)})^2 \right] \geq \rho_{\text{CNN}} (\|\mathbf{w}^{\star}\|_2). \tag{B.47}$$

Hence,

$$\nabla^2 f(\mathbf{w}^*) \succeq \frac{4}{K} \cdot \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2) \cdot \mathbf{I}. \quad (\text{B.48})$$

Moreover, we apply Cauchy-Schwarz inequality and upper bound the Hessian as

$$\begin{aligned} \nabla^2 f(\mathbf{w}^*) &\leq \left(\max_{\|\mathbf{u}\|_2=1} \mathbf{u}^\top \nabla^2 f(\mathbf{w}^*) \mathbf{u} \right) \cdot \mathbf{I} \\ &\leq \max_{\|\mathbf{u}\|_2=1} \mathbb{E} \left[\frac{\sum_{j=1}^K \left(\frac{1}{K} \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(j)}) \right)^2 \cdot \sum_{j=1}^K (\mathbf{u}^\top \mathbf{x}^{(j)})^2}{H(\mathbf{w}^*) (1 - H(\mathbf{w}^*))} \right] \cdot \mathbf{I}. \end{aligned} \quad (\text{B.49})$$

Using (B.40), i.e.,

$$\frac{\frac{1}{K^2} \sum_{j=1}^K \phi'(\mathbf{w}^{*\top} \mathbf{x}^{(j)})^2}{H(\mathbf{w}^*) (1 - H(\mathbf{w}^*))} \leq \frac{1}{4}, \quad (\text{B.50})$$

we upper bound the right-hand side of (B.49) as

$$\nabla^2 f(\mathbf{w}^*) \preceq \max_{\|\mathbf{u}\|_2=1} \mathbb{E} \left[\frac{1}{4} \sum_{j=1}^K (\mathbf{u}^\top \mathbf{x}^{(j)})^2 \right] \cdot \mathbf{I} = \frac{K}{4} \cdot \mathbf{I}. \quad (\text{B.51})$$

Together with the lower bound, we now conclude that

$$\frac{4}{K} \cdot \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2) \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{w}^*) \preceq \frac{K}{4} \cdot \mathbf{I}. \quad (\text{B.52})$$

And following from (B.6) in Lemma 5, we have

$$\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| \leq C \cdot K \cdot \|\mathbf{w} - \mathbf{w}^*\|_2. \quad (\text{B.53})$$

Thus if $\|\mathbf{w} - \mathbf{w}^*\| \leq \min \left\{ 0.7, C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K^2} \right\}$, we have

$$C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{w}) \preceq C \cdot K \cdot \mathbf{I}. \quad (\text{B.54})$$

B.4.3 Proof of Lemma 7

We apply a covering type of argument to show that the Hessian of the empirical risk function concentrates around the Hessian of the population risk function uniformly, and the argument applies to both the loss associated with FCN and CNN. We first take the FCN case as an example and then we provide the necessary modifications for the proof of the CNN case.

- **The FCN case:** We adapt the analysis in [50] to our setting. Let N_ϵ be the ϵ -covering number of the Euclidean ball $\mathbb{B}(\mathbf{W}^*, r)$. Here, we omit the subscript FCN of r for simplicity. It is known that $\log N_\epsilon \leq dK \log(3r/\epsilon)$ [85]. Let $\mathcal{W}_\epsilon = \{\mathbf{W}_1, \dots, \mathbf{W}_{N_\epsilon}\}$ be the ϵ -cover set with N_ϵ elements. For any $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)$, let $j(\mathbf{W}) = \operatorname{argmin}_{j \in [N_\epsilon]} \|\mathbf{W} - \mathbf{W}_{j(\mathbf{W})}\|_F \leq \epsilon$ for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)$.

For any $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)$, we have

$$\begin{aligned} \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| &\leq \frac{1}{n} \left\| \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i)] \right\| \\ &\quad + \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i) - \mathbb{E} [\nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x})] \right\| \\ &\quad + \|\mathbb{E} [\nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x})] - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})]\|. \end{aligned}$$

Hence, we have

$$\mathbb{P} \left(\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| \geq t \right) \leq \mathbb{P}(A_t) + \mathbb{P}(B_t) + \mathbb{P}(C_t),$$

where the events A_t , B_t and C_t are defined as

$$\begin{aligned} A_t &= \left\{ \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{W}^*, r)} \frac{1}{n} \left\| \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x}_i)] \right\| \geq \frac{t}{3} \right\}, \\ B_t &= \left\{ \sup_{\mathbf{w} \in \mathcal{W}_\epsilon} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\| \geq \frac{t}{3} \right\}, \\ C_t &= \left\{ \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{W}^*, r)} \left\| \mathbb{E} [\nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x})] - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\| \geq \frac{t}{3} \right\}. \end{aligned}$$

In the sequel, we bound the terms $\mathbb{P}(A_t)$, $\mathbb{P}(B_t)$, and $\mathbb{P}(C_t)$, separately.

1. **Upper bound on $\mathbb{P}(B_t)$.** Before continuing, we state a useful technical lemma, whose proof can be found in [50].

Lemma 10. *Let $\mathbf{M} \in \mathbb{R}^{d \times d}$ be a symmetric $d \times d$ matrix and V_ϵ be an ϵ -cover of unit-Euclidean-norm ball $\mathbb{B}(\mathbf{0}, 1)$, then*

$$\|\mathbf{M}\| \leq \frac{1}{1 - 2\epsilon} \sup_{\mathbf{v} \in V_\epsilon} |\langle \mathbf{v}, \mathbf{M}\mathbf{v} \rangle|. \quad (\text{B.55})$$

Let $V_{\frac{1}{4}}$ be a $(\frac{1}{4})$ -cover of the ball $\mathbb{B}(\mathbf{0}, 1) = \{\mathbf{W} \in \mathbb{R}^{d \times K} : \|\mathbf{W}\|_{\text{F}} = 1\}$, where $\log |V_{\frac{1}{4}}| \leq dK \log 12$. Following from Lemma 10, we have

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\| \\ & \leq 2 \sup_{\mathbf{v} \in V_{\frac{1}{4}}} \left| \left\langle \mathbf{v}, \left(\frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right) \mathbf{v} \right\rangle \right|. \end{aligned}$$

Taking the union bound over \mathcal{W}_ϵ and $V_{\frac{1}{4}}$ yields

$$\begin{aligned} \mathbb{P}(B_t) &\leq \mathbb{P} \left(\sup_{\mathbf{w} \in \mathcal{W}_\epsilon, \mathbf{v} \in V_{\frac{1}{4}}} \left| \frac{1}{n} \sum_{i=1}^n G_i \right| \geq \frac{t}{6} \right) \\ &\leq e^{dK(\log \frac{3r}{\epsilon} + \log 12)} \sup_{\mathbf{w} \in \mathcal{W}_\epsilon, \mathbf{v} \in V_{\frac{1}{4}}} \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n G_i \right| \geq \frac{t}{6} \right), \quad (\text{B.56}) \end{aligned}$$

where $G_i = \langle \mathbf{v}, (\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})]) \mathbf{v} \rangle$ and $\mathbb{E}[G_i] = 0$. Let $\mathbf{a} = [\mathbf{a}_1^\top, \dots, \mathbf{a}_K^\top] \in \mathbb{R}^{dK}$. Then we can show that $\|G_i\|_{\psi_1}$ is upper bounded, which we summarize as follows, and whose proof is given in Appendix B.4.6.

Lemma 11. *Suppose the loss is associated with FCN. There exists some constant C such that*

$$\|G_i\|_{\psi_1} \leq C \equiv \tau^2.$$

Applying the Bernstein inequality for sub-exponential random variables [50, Theorem 9] to (B.56), we have that for fixed $\mathbf{W} \in \mathcal{W}_\epsilon$, $\mathbf{v} \in V_{\frac{1}{4}}$,

$$\begin{aligned} & \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{v}, (\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})]) \mathbf{v} \rangle \right| \geq \frac{t}{6} \right) \\ & \leq 2 \exp \left(-c \cdot n \cdot \min \left(\frac{t^2}{\tau^4}, \frac{t}{\tau^2} \right) \right), \end{aligned} \quad (\text{B.57})$$

for some universal constant c . As a result, $\mathbb{P}(B_t)$ is upper bounded by

$$2 \exp \left(-c \cdot n \cdot \min \left(\frac{t^2}{\tau^4}, \frac{t}{\tau^2} \right) + dK \log \frac{3r}{\epsilon} + dK \log 12 \right).$$

Thus as long as

$$t > C \cdot \max \left\{ \sqrt{\frac{\tau^4 (dK \log \frac{36r}{\epsilon} + \log \frac{4}{\delta})}{n}}, \frac{\tau^2 (dK \log \frac{36r}{\epsilon} + \log \frac{4}{\delta})}{n} \right\} \quad (\text{B.58})$$

for some large enough constant C , we have $\mathbb{P}(B_t) \leq \frac{\delta}{2}$.

2. **Upper bound on $\mathbb{P}(A_t)$ and $\mathbb{P}(C_t)$.** These two events will be bounded in a similar way. We first present the following useful Lemma, whose proof is provided in Appendix B.4.8

Lemma 12. *Suppose the loss is associated with FCN. There exists some constant C such that*

$$\mathbb{E} \left[\sup_{\mathbf{W} \neq \mathbf{W}' \in \mathbb{B}(\mathbf{W}^*, r)} \frac{\|\nabla^2 \ell(\mathbf{W}, \mathbf{x}) - \nabla^2 \ell(\mathbf{W}', \mathbf{x})\|}{\|\mathbf{W} - \mathbf{W}'\|_F} \right] \leq C \cdot d\sqrt{K}. \quad (\text{B.59})$$

Consider the event C_t first. We derive

$$\begin{aligned}
& \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \left\| \mathbb{E} [\nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x})] - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\| \\
& \leq \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \frac{\left\| \mathbb{E} [\nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x})] - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\|}{\|\mathbf{W} - \mathbf{W}_{j(\mathbf{w})}\|_{\text{F}}} \\
& \quad \cdot \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \|\mathbf{W} - \mathbf{W}_{j(\mathbf{w})}\|_{\text{F}} \\
& \leq C \cdot d\sqrt{K} \cdot \epsilon.
\end{aligned} \tag{B.60}$$

Therefore, C_t holds as long as

$$t \geq C \cdot d\sqrt{K} \cdot \epsilon. \tag{B.61}$$

We can bound the event A_t as below.

$$\begin{aligned}
& \mathbb{P} \left(\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \frac{1}{n} \left\| \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x}_i)] \right\| \geq \frac{t}{3} \right) \\
& \leq \frac{3}{t} \mathbb{E} \left[\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \left\| \frac{1}{n} \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x}_i)] \right\| \right] \\
& \leq \frac{3}{t} \mathbb{E} \left[\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \left\| \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x}_i) \right\| \right] \\
& \leq \frac{3}{t} \mathbb{E} \left[\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \frac{\left\| \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{w})}; \mathbf{x}_i) \right\|}{\|\mathbf{W} - \mathbf{W}_{j(\mathbf{w})}\|_{\text{F}}} \right] \\
& \quad \cdot \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \|\mathbf{W} - \mathbf{W}_{j(\mathbf{w})}\|_{\text{F}} \\
& \leq \frac{C \cdot d\sqrt{K} \cdot \epsilon}{t}
\end{aligned} \tag{B.62}$$

where (B.62) follows from the Markov inequality. Thus, taking

$$t \geq \frac{6\epsilon \cdot C \cdot d\sqrt{K}}{\delta} \tag{B.64}$$

ensures that $\mathbb{P}(A_t) \leq \frac{\delta}{2}$.

3. **Final step.** Let $\epsilon = \frac{\delta\tau^2}{C \cdot d\sqrt{K} \cdot ndK}$ and $\delta = d^{-10}$. Plugging ϵ and δ into (B.58)

we need

$$t > \tau^2 \cdot \max \left\{ \frac{1}{ndK}, C \cdot \sqrt{\frac{(dK \log(36rnd^{11}K) + \log \frac{4}{\delta})}{n}}, \frac{(dK \log(36rnd^{11}K) + \log \frac{4}{\delta})}{n} \right\}.$$

The middle term can be bounded as

$$\frac{dK \log(36rnd^{11}K) + 10 \log d}{n} \leq \frac{dK \log n}{n} + \frac{dK \log 36r}{n} + \frac{11dK \log dK}{n} + \frac{10 \log d}{n}.$$

If $n \geq C \cdot dK \log(dK)$ for some large enough constant C , the first term $dK \log n$ dominates and is on the order of $dK \log(dK)$. Moreover, it decreases as n increases when $n \geq 3$. Thus we can set

$$t \geq \tau^2 \sqrt{\frac{(dK \log(36rnd^{11}K) + \log \frac{4}{\delta})}{n}} \quad (\text{B.65})$$

which holds as $t \geq C' \cdot \tau^2 \sqrt{\frac{dK \log n}{n}}$ for some constant C' . By setting $t := C\tau^2 \sqrt{\frac{dK \log n}{n}}$ for sufficiently large C , as long as $n \geq C' \cdot dK \log dK$,

$$\mathbb{P} \left(\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \|\nabla^2 f_n(\mathbf{w}) - \nabla^2 f(\mathbf{w})\| \geq C\tau^2 \sqrt{\frac{dK \log n}{n}} \right) \leq d^{-10}.$$

- **The CNN case:** If the loss is associated with CNN, we redefine G_i as $G_i = \langle \mathbf{v}, (\nabla^2 \ell(\mathbf{w}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{w}; \mathbf{x})]) \mathbf{v} \rangle$ and we show the following Lemmas whose proof is given in Appendix B.4.7 and Appendix B.4.9.

Lemma 13. *Suppose the loss is associated CNN. There exists some constant C such that*

$$\|G_i\|_{\psi_1} \leq C \cdot K^2 := \tau^2. \quad (\text{B.66})$$

Lemma 14. *Suppose the loss is associated with CNN. There exists some constant C such that*

$$\mathbb{E} \left[\sup_{\mathbf{W} \neq \mathbf{W}' \in \mathbb{B}(\mathbf{W}^*, r)} \frac{\|\nabla^2 \ell(\mathbf{W}, \mathbf{x}) - \nabla^2 \ell(\mathbf{W}', \mathbf{x})\|}{\|\mathbf{W} - \mathbf{W}'\|_F} \right] \leq C \cdot d\sqrt{K}. \quad (\text{B.67})$$

Following argument similar to the proof of Lemma 7, we can obtain the following concentration inequality:

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \|\nabla^2 f_n(\mathbf{w}) - \nabla^2 f(\mathbf{w})\| \leq C \cdot K^2 \sqrt{\frac{\frac{d}{K} \cdot \log n}{n}}, \quad (\text{B.68})$$

holds with probability at least $1 - d^{-10}$, as long as the sample complexity $n \geq C \cdot \frac{d}{K} \log\left(\frac{d}{K}\right)$.

B.4.4 Proof of Lemma 8

In order to proceed we need the following Lemma 15 whose proof is given in Appendix B.4.10.

Lemma 15. *Suppose the loss is associated with FCN. Let \mathbf{u} be a fixed unit norm vector $\mathbf{u} = [\mathbf{u}_1^\top, \dots, \mathbf{u}_K^\top] \in \mathbb{R}^{dK}$ with $\|\mathbf{u}\|_2 = 1$. Then we have*

$$\|\mathbf{u}^\top \nabla \ell(\mathbf{W}; \mathbf{x})\|_{\psi_2} \leq \sqrt{K}.$$

Suppose the loss is associated with CNN. Let \mathbf{u} be a fixed unit norm vector $\mathbf{u} \in \mathbb{R}^m$ with $\|\mathbf{u}\|_2 = 1$. Then

$$\|\langle \mathbf{u}, \nabla \ell(\mathbf{w}) \rangle\|_{\psi_2} \leq C \cdot K.$$

Following argument (details omitted) similar to the proof of Lemma 7, and applies Lemma 15, for the loss associated with FCN, we can get the following concentration

inequality

$$\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})} \|\nabla f_n(\mathbf{W}) - \nabla f(\mathbf{W})\|_2 \leq C \cdot \sqrt{\frac{d\sqrt{K} \log n}{n}} \quad (\text{B.69})$$

with probability at least $1 - d^{-10}$, as long as the sample size $n \geq C \cdot dK \log(dK)$. For the loss associated with CNN, we obtain

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})} \|\nabla f_n(\mathbf{w}) - \nabla f(\mathbf{w})\| \leq C \cdot \sqrt{K} \sqrt{\frac{\frac{d}{K} \log n}{n}} = C \cdot \sqrt{\frac{d \log n}{n}}, \quad (\text{B.70})$$

with probability at least $1 - d^{-10}$ as long as $n \geq C \cdot \frac{d}{K} \log\left(\frac{d}{K}\right)$.

B.4.5 Proof of Lemma 9

We take the first term in (B.26) as an example, since the second term follows exactly in the same way. We first derive

$$\mathbb{E} \left[\left(\frac{1}{K} \sum_{i=1}^K \phi(\mathbf{w}_i^\top \mathbf{x}) \right)^{-z} \right] \leq \mathbb{E} \left[\frac{1}{K} \sum_{i=1}^K (\phi(\mathbf{w}_i^\top \mathbf{x}))^{-z} \right], \quad (\text{B.71})$$

which follows from the fact that $f(x) = x^{-z}$ is convex for $x > 0$ and $z \geq 1$. Further since $\frac{1}{\phi(x)} = 1 + e^{-x}$, and $g = \mathbf{w}_i^\top \mathbf{x} \sim \mathcal{N}(0, \sigma_i^2 = \|\mathbf{w}_i\|_2^2)$, we can exactly calculate the summands in the above equation as follows:

$$\mathbb{E} [\phi(g)^{-z}] = \mathbb{E} \left[\sum_{l=0}^z \binom{z}{l} e^{-lg} \right] = \sum_{l=0}^z \binom{z}{l} e^{\left(\frac{\sigma_i^2 l^2}{2}\right)},$$

where we use the fact that g is a Gaussian random variable. Hence, we conclude that for $t = \max(\|\mathbf{w}_1\|_2, \dots, \|\mathbf{w}_K\|_2)$ and $p \geq 1$,

$$\mathbb{E} \left[\left(\frac{1}{\frac{1}{K} \sum_{i=1}^K \phi(\mathbf{w}_i^\top \mathbf{x})} \right)^z \right] \leq C \cdot e^{t^2}, \quad (\text{B.72})$$

holds for some constant C depending on z .

B.4.6 Proof of Lemma 11

The sub-exponential norm of G_i can be bounded as

$$\|G_i\|_{\psi_1} \leq \|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{W}; z) \mathbf{u} \rangle\|_{\psi_1} + \|\nabla^2 f(\mathbf{W}; z)\|,$$

where $\|\nabla^2 f(\mathbf{W}; z)\|$ is upper bounded by C due to Lemma 6. Denote the (j, l) -th block of $\nabla^2 \ell(\mathbf{W}; z)$ as $\xi_{j,l} \cdot \mathbf{x} \mathbf{x}^\top$. We can derive

$$\begin{aligned} \|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{W}; z) \mathbf{u} \rangle\|_{\psi_1} &\leq \sum_{j=1}^K \sum_{l=1}^K \|\xi_{j,l} \cdot \mathbf{u}_j^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}_l\|_{\psi_1} \\ &\leq \sum_{j=1}^K \sum_{l=1}^K \sup_{t \geq 1} t^{-1} \left(\mathbb{E} |\xi_{j,l} \cdot \mathbf{u}_j^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}_l|^t \right)^{\frac{1}{t}}. \end{aligned} \quad (\text{B.73})$$

Next we show that $\xi_{j,l}$ is upper bounded by some constant for all j and l .

- For $j \neq l$,

$$\begin{aligned} |\xi_{j,l}| &= \left| \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x}) \cdot (H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W}))}{(H(\mathbf{W})(1 - H(\mathbf{W})))^2} \right| \\ &= \begin{cases} \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{(1 - H(\mathbf{W}))^2} & y = 0 \\ \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{H(\mathbf{W})^2} & y = 1 \end{cases}. \end{aligned} \quad (\text{B.74})$$

Moreover,

$$\frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{(1 - H(\mathbf{W}))^2} \leq \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{(1 - \phi(\mathbf{w}_j^\top \mathbf{x})) (1 - \phi(\mathbf{w}_l^\top \mathbf{x}))} \leq \phi(\mathbf{w}_j^\top \mathbf{x}) \phi(\mathbf{w}_l^\top \mathbf{x}) \leq 1, \quad (\text{B.75})$$

where the first inequality holds due to the following fact,

$$(1 - H(\mathbf{W}))^2 = \left(1 - \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_j^\top \mathbf{x}) \right)^2 \geq \frac{1}{K^2} (1 - \phi(\mathbf{w}_j^\top \mathbf{x})) (1 - \phi(\mathbf{w}_l^\top \mathbf{x})),$$

the second inequality follows because $\phi(x)(1 - \phi(x)) = \phi'(x)$. Similarly, we can show that

$$\frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{H(\mathbf{W})^2} \leq 1. \quad (\text{B.76})$$

Thus for $j \neq l$, $|\xi_{j,l}| \leq 1$ holds.

- For $j = l$,

$$|\xi_{j,j}| \leq \left| \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_j^\top \mathbf{x}) \cdot (H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W}))}{(H(\mathbf{W})(1 - H(\mathbf{W})))^2} \right| + \left| \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x}) (y - H(\mathbf{W}))}{H(\mathbf{W})(1 - H(\mathbf{W}))} \right|.$$

For the second term in the above equation, we have

$$\left| \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x}) (y - H(\mathbf{W}))}{H(\mathbf{W})(1 - H(\mathbf{W}))} \right| = \begin{cases} \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{(1 - H(\mathbf{W}))} \leq 1 & y = 0 \\ \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})} \leq 1 & y = 1 \end{cases},$$

which follows from the fact that the second derivative is

$$\phi''(x) = \phi(x)(1 - \phi(x))(1 - 2\phi(x)),$$

the absolute value of which can be upper bounded by $\phi(x)$ or $1 - \phi(x)$.

Hence,

$$\begin{aligned} & \|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{W}; z) \mathbf{u} \rangle\|_{\psi_1} \\ & \leq C \cdot \sum_{j=1}^K \sum_{l=1}^K \sup_{t \geq 1} t^{-1} \left(\sqrt{\mathbb{E}[(\mathbf{u}_j^\top \mathbf{x})^{2t}]} \cdot \sqrt{\mathbb{E}[(\mathbf{u}_l^\top \mathbf{x})^{2t}]} \right)^{\frac{1}{t}} \\ & \leq C \cdot \sum_{j=1}^K \sum_{l=1}^K \|\mathbf{u}_j\|_2 \|\mathbf{u}_l\|_2 \cdot \sup_{t \geq 1} t^{-1} ((2t - 1)!!)^{\frac{1}{t}} \\ & \leq C := \tau^2, \end{aligned} \tag{B.77}$$

where the last inequality holds because

$$\begin{aligned} \sup_{t \geq 1} t^{-1} ((2t - 1)!!)^{\frac{1}{t}} & \leq \sup_{t \geq 1} t^{-1} ((2t)^t)^{\frac{1}{t}} \leq 2, \\ \sum_{j=1}^K \sum_{l=1}^K \|\mathbf{u}_j\|_2 \|\mathbf{u}_l\|_2 & \leq \sum_{j=1}^K \sum_{l=1}^K \frac{\|\mathbf{u}_j\|_2^2 + \|\mathbf{u}_l\|_2^2}{2} = \frac{1}{2}. \end{aligned} \tag{B.78}$$

Thus, we conclude

$$\|G_i\|_{\psi_1} \leq C := \tau^2.$$

B.4.7 Proof of Lemma 13

Again the sub-exponential norm of G_i can be bounded as

$$\|G_i\|_{\psi_1} \leq \|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}; z) \mathbf{u} \rangle\|_{\psi_1} + \|\nabla^2 f(\mathbf{w}; z)\|,$$

where $\|\nabla^2 f(\mathbf{W}; z)\|$ is upper bounded by $C \cdot K$ due to Lemma 6. Applying the triangle inequality, the sub-exponential norm of $\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}) \mathbf{u} \rangle$ can be bounded as

$$\|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}) \mathbf{u} \rangle\|_{\psi_1} \leq \sum_{j \neq l} \|g_{j,l}(\mathbf{w}) \mathbf{u}^\top \mathbf{x}^{(j)} \mathbf{u}^\top \mathbf{x}^{(l)}\|_{\psi_1} + \sum_{j=l} \|g_{j,l}(\mathbf{w}) \mathbf{u}^\top \mathbf{x}^{(j)} \mathbf{u}^\top \mathbf{x}^{(l)}\|_{\psi_1}. \quad (\text{B.79})$$

Hence, we have

$$\begin{aligned} & \left| \frac{1}{K^2} \frac{H(\mathbf{w})^2 + y - 2y \cdot H(\mathbf{w})}{(H(\mathbf{w})(1 - H(\mathbf{w})))^2} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \right| \\ &= \begin{cases} \frac{1}{K^2} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)})}{H(\mathbf{w})^2} \leq 1 & y = 1 \\ \frac{1}{K^2} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)})}{(1 - H(\mathbf{w}))^2} \leq 1 & y = 0 \end{cases}, \end{aligned}$$

and

$$\left| \frac{1}{K} \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \right| = \begin{cases} \frac{1}{K} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{H(\mathbf{w})} \leq 1 & y = 1 \\ \frac{1}{K} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{1 - H(\mathbf{w})} \leq 1 & y = 0 \end{cases}.$$

Plugging it back to (B.79), we obtain

$$\|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}) \mathbf{u} \rangle\|_{\psi_1} \leq \sum_{j \neq l} \|(\mathbf{u}^\top \mathbf{x}^{(j)}) (\mathbf{u}^\top \mathbf{x}^{(l)})\|_{\psi_1} + \sum_{j=1}^K \|(\mathbf{u}^\top \mathbf{x}^{(j)})^2\|_{\psi_1} \leq C \cdot K^2. \quad (\text{B.80})$$

B.4.8 Proof of Lemma 12

As noted before, we can write the (j, l) -th block of $\nabla^2 \ell(\mathbf{W}; z)$ as $\xi_{j,l}(\mathbf{W}) \mathbf{x} \mathbf{x}^\top$.

Then we can obtain the following bound,

$$\|\nabla^2 \ell(\mathbf{W}; z) - \nabla^2 \ell(\mathbf{W}'; z)\| \leq \sum_{j=1}^K \sum_{l=1}^K |\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}')| \cdot \|\mathbf{x} \mathbf{x}^\top\|. \quad (\text{B.81})$$

Using the same method as shown in the proof of Lemma 5, we can upper bound $|\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}')|$ as

$$|\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}')| \leq \left(\max_k |T_{j,l,k}| \right) \cdot \|\mathbf{x}\|_2 \cdot \sqrt{K} \cdot \|\mathbf{W} - \mathbf{W}'\|_F,$$

where following from (B.23),

$$|T_{j,l,k}| \leq \max \left\{ \frac{2}{K^3} \frac{1}{H(\mathbf{W})^3}, \frac{1}{K^2} \frac{1}{H(\mathbf{W})^2}, \frac{2}{K^3} \frac{1}{(1-H(\mathbf{W}))^3}, \frac{1}{K^2} \frac{1}{(1-H(\mathbf{W}))^2} \right\}. \quad (\text{B.82})$$

And thus, if $\|\mathbf{W} - \mathbf{W}'\|_F \leq 0.7$ we have

$$\begin{aligned} \mathbb{E} \left[\sup_{\mathbf{W} \neq \mathbf{W}'} \frac{\|\nabla^2 \ell(\mathbf{W}) - \nabla^2 \ell(\mathbf{W}')\|}{\|\mathbf{W} - \mathbf{W}'\|_F} \right] &\leq \sqrt{K} \cdot K^2 \cdot \mathbb{E} \left[\left(\max_{j,l,k} |T_{j,l,k}| \right) \cdot \|\mathbf{x}\|_2 \cdot \|\mathbf{x}\mathbf{x}^\top\| \right] \\ &\leq C \cdot d\sqrt{K}. \end{aligned} \quad (\text{B.83})$$

Thus we only need to set $J^\star \geq C \cdot d\sqrt{K}$ for some large enough C .

B.4.9 Proof of Lemma 14

Following from (B.4) we can write

$$\|\nabla^2 \ell(\mathbf{w}) - \nabla^2 \ell(\mathbf{w}')\| \leq \sum_{j=1}^K \sum_{l=1}^K |g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}')| \cdot \|\mathbf{x}^{(j)} \mathbf{x}^{(l)\top}\|. \quad (\text{B.84})$$

Similarly, the analysis in the proof of Lemma 5 implies that

$$|g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}')| \leq \left(\max_k |S_{j,l,k}| \right) \cdot \sqrt{K} \|\mathbf{x}\|_2 \cdot \|\mathbf{w} - \mathbf{w}'\|_2, \quad (\text{B.85})$$

where we upper-bound $S_{j,l,k}$ in (B.34) as

$$|S_{j,l,k}| \leq \begin{cases} \max \left\{ \frac{1}{K^2} \frac{1}{(1-H(\mathbf{w}))^3}, \frac{1}{K^2} \frac{1}{(H(\mathbf{w}))^3} \right\} & j \neq l \\ \max \left\{ \frac{1}{K} \frac{1}{(1-H(\mathbf{w}))^2}, \frac{1}{K} \frac{1}{(H(\mathbf{w}))^2} \right\} & j = l \end{cases}. \quad (\text{B.86})$$

Hence, if $\|\mathbf{w} - \mathbf{w}'\|_2 \leq 0.7$, we have

$$\begin{aligned} \mathbb{E} \left[\sup_{\mathbf{w} \neq \mathbf{w}'} \frac{\|\nabla^2 \ell(\mathbf{w}) - \nabla^2 \ell(\mathbf{w}')\|}{\|\mathbf{w} - \mathbf{w}'\|_F} \right] &\leq \sqrt{K} \cdot \sum_{j=1}^K \sum_{l=1}^K \mathbb{E} \left[\left(\max_k |S_{j,l,k}| \right) \cdot \|\mathbf{x}\|_2 \cdot \|\mathbf{x}^{(j)} \mathbf{x}^{(l)\top}\| \right] \\ &\leq C \cdot d\sqrt{K}. \end{aligned} \quad (\text{B.87})$$

Thus, in this case we can set $J^\star \geq C \cdot d\sqrt{K}$ as well.

B.4.10 Proof of Lemma 15

- **The FCN case:** Following from (B.1), we have

$$\langle \nabla \ell(\mathbf{W}), \mathbf{u} \rangle = \frac{1}{K} \sum_{j=1}^K \left(\frac{(y - H(\mathbf{W})) \cdot \phi'(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})(1 - H(\mathbf{W}))} \right) (\mathbf{u}_j^\top \mathbf{x}),$$

and by definition, we can upper-bound the sub-Gaussian norm as

$$\begin{aligned} & \|\langle \nabla \ell(\mathbf{W}), \mathbf{u} \rangle\|_{\psi_2} \\ & \leq \begin{cases} \frac{1}{K} \sum_{j=1}^K \left\| \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})}{(1 - \frac{1}{K} \sum_{l=1}^K \phi(\mathbf{w}_l^\top \mathbf{x}))} \mathbf{u}_j^\top \mathbf{x} \right\|_{\psi_2} & y = 0 \\ \frac{1}{K} \sum_{j=1}^K \left\| \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})}{\frac{1}{K} \sum_{l=1}^K \phi(\mathbf{w}_l^\top \mathbf{x})} \mathbf{u}_j^\top \mathbf{x} \right\|_{\psi_2} & y = 1 \end{cases}. \end{aligned}$$

Thus we conclude that

$$\|\langle \nabla \ell(\mathbf{W}), \mathbf{u} \rangle\|_{\psi_2} \leq \sum_{j=1}^K \|\mathbf{u}_j\|_2 \leq \sqrt{K}, \quad (\text{B.88})$$

and the directional gradient is \sqrt{K} -sub-Gaussian.

- **The CNN case:** Following from (B.3), we have

$$\langle \nabla \ell(\mathbf{w}), \mathbf{u} \rangle = - \sum_{j=1}^K \frac{1}{K} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \cdot (\mathbf{u}^\top \mathbf{x}^{(j)}),$$

where

$$\left| \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \right| = \begin{cases} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{\sum_{j=1}^K \frac{1}{K} \phi(\mathbf{w}^\top \mathbf{x}^{(j)})} \leq K & y = 1 \\ \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{\sum_{j=1}^K \frac{1}{K} (1 - \phi(\mathbf{w}^\top \mathbf{x}^{(j)}))} \leq K & y = 0 \end{cases}.$$

Then the sub-Gaussian norm of $\langle \nabla \ell(\mathbf{w}), \mathbf{u} \rangle$ is upper bounded as

$$\|\langle \nabla \ell(\mathbf{w}), \mathbf{u} \rangle\|_{\psi_2} \leq K \cdot \frac{1}{K} \sum_{j=1}^K \|\mathbf{u}^\top \mathbf{x}^{(j)}\|_{\psi_2} \leq C \cdot K. \quad (\text{B.89})$$

Hence, the directional gradient is K -sub-Gaussian.

B.5 Proof of Theorem 5

We first define a product $\tilde{\otimes}$ as follows. If $\mathbf{v} \in \mathbb{R}^d$ is a vector and \mathbf{I} is the identity matrix, then $\mathbf{v} \tilde{\otimes} \mathbf{I} = \sum_{j=1}^d [\mathbf{v} \otimes \mathbf{e}_j \otimes \mathbf{e}_j + \mathbf{e}_j \otimes \mathbf{v} \otimes \mathbf{e}_j + \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}]$. If \mathbf{M} is a symmetric rank- r matrix factorized as $\mathbf{M} = \sum_{i=1}^r \mathbf{s}_i \mathbf{v}_i \mathbf{v}_i^\top$ and \mathbf{I} is the identity matrix, then

$$\mathbf{M} \tilde{\otimes} \mathbf{I} = \sum_{i=1}^r \mathbf{s}_i \sum_{j=1}^d \sum_{l=1}^6 \mathbf{A}_{l,i,j}, \quad (\text{B.90})$$

where $\mathbf{A}_{1,i,j} = \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_j$, $\mathbf{A}_{2,i,j} = \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{e}_j$, $\mathbf{A}_{3,i,j} = \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{e}_j$, $\mathbf{A}_{4,i,j} = \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}_i$, $\mathbf{A}_{5,i,j} = \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{v}_i$ and $\mathbf{A}_{6,i,j} = \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{v}_i$.

And We further define a tensor operation as follows. For a tensor $\mathbf{T} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ and three matrices $\mathbf{A} \in \mathbb{R}^{n_1 \times d_1}$, $\mathbf{B} \in \mathbb{R}^{n_2 \times d_2}$, $\mathbf{C} \in \mathbb{R}^{n_3 \times d_3}$, the (i, j, k) -th entry of the tensor $\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C})$ is given by

$$\sum_{i'}^{n_1} \sum_{j'}^{n_2} \sum_{k'}^{n_3} \mathbf{T}_{i',j',k'} \mathbf{A}_{i',i} \mathbf{B}_{j',j} \mathbf{C}_{k',k}. \quad (\text{B.91})$$

The proof contains two parts. Part (a) proves that the estimation of the direction of \mathbf{W}^* is sufficiently accurate, which follows from the arguments similar to those in [7] and is only briefly summarized below. Part (b) is different, where we do not require the homogeneous condition for the activation function, and instead, our proof is based on a mild condition in Assumption 2. We detail our proof in part (b).

(a) In order to estimate the direction of each \mathbf{w}_i for $i = 1, \dots, K$, [7] showed that for the regression problem, if the sample size $n \geq d \text{poly}(K, \kappa, \zeta, \log d)$, where $\zeta > 1$ is any constant, then

$$\|\overline{\mathbf{w}}_i^* - s_i \mathbf{V} \hat{\mathbf{u}}_i\| \leq \epsilon \text{poly}(K, \kappa) \quad (\text{B.92})$$

holds with probability at least $1 - d^{-\Omega(\zeta)}$. Such a result also holds for the classification problem with only slight difference in the proof as we describe as follows. The main

idea of the proof is to bound the estimation error of \mathbf{P}_2 and \mathbf{R}_3 via Bernstein inequality. For the regression problem, Bernstein inequality was applied to terms associated with each neuron individually, and the bounds were then put together via the triangle inequality in [7]. However, for the classification problem here, we apply Bernstein inequality to the terms associated with all neurons together. Another difference is that the label y_i of the classification model is bounded by nature, whereas the output y_i in the regression model needs to be upper-bounded via homogeneously bounded conditions of the activation function. A reader can refer to [7] for the details of the proof for this part.

(b) In order to estimate $\|\mathbf{w}_i\|$ for $i = 1, \dots, K$, we provide a different proof from [7], which does not require the homogeneous condition on the activation function, but assumes a more relaxed condition in Assumption 2.

We define a quantity Q_1 as follows:

$$Q_1 = \mathbf{M}_{l_1}(\mathbf{I}, \underbrace{\boldsymbol{\alpha}, \dots, \boldsymbol{\alpha}}_{(l_1-1)}), \quad (\text{B.93})$$

where l_1 is the first non-zero index such that $\mathbf{M}_{l_1} \neq 0$. For example, if $l_1 = 3$, then Q_1 takes the following form

$$Q_1 = \mathbf{M}_3(\mathbf{I}, \boldsymbol{\alpha}, \boldsymbol{\alpha}) = \frac{1}{K} \sum_{i=1}^K m_{3,i}(\|\mathbf{w}_i^*\|) (\boldsymbol{\alpha}^\top \bar{\mathbf{w}}_i^*)^2 \bar{\mathbf{w}}_i^*, \quad (\text{B.94})$$

where $\bar{\mathbf{w}} = \mathbf{w}/\|\mathbf{w}\|$ and by definition

$$m_{3,i}(\|\mathbf{w}_i^*\|) = \mathbb{E} [\phi(\|\mathbf{w}_i^*\| \cdot z) z^3] - 3\mathbb{E} [\phi(\|\mathbf{w}_i^*\| \cdot z) z]. \quad (\text{B.95})$$

Clearly, Q_1 has information of $\|\mathbf{w}_i^*\|$, which can be estimated by solving the following optimization problem:

$$\beta^* = \operatorname{argmin}_{\beta \in \mathbb{R}^K} \left\| \frac{1}{K} \sum_{i=1}^K \beta_i s_i \bar{\mathbf{w}}_i^* - Q_1 \right\|, \quad (\text{B.96})$$

where each entry of the solution takes the form

$$\beta_i^* = s_i^3 m_{3,i}(\|\mathbf{w}_i^*\|) (\boldsymbol{\alpha}^T s_i \overline{\mathbf{w}_i^*})^2. \quad (\text{B.97})$$

In the initialization, we substitute \widehat{Q}_1 (estimated from training data) for Q_1 , $\mathbf{V}\hat{u}_i$ (estimated in part (a)) for $s_i \overline{\mathbf{w}_i^*}$ into (B.96), and obtain an estimate $\hat{\beta}$ of β^* . We then substitute $\hat{\beta}$ for β^* and $\mathbf{V}\hat{u}_i$ for $s_i \overline{\mathbf{w}_i^*}$ into (B.97) to obtain an estimate \hat{a}_i of $\|\mathbf{w}_i^*\|$ via the following equation

$$\hat{\beta}_i = s_i^3 m_{3,i}(\hat{a}_i) (\boldsymbol{\alpha}^T \mathbf{V}\hat{u}_i)^2. \quad (\text{B.98})$$

Furthermore, since $m_{l_1,i}(x)$ has fixed sign for $x > 0$ and for $l_1 \geq 1$, s_i can be estimated correctly from the sign of $\hat{\beta}_i$ for $i = 1, \dots, K$.

For notational simplicity, let $\beta_{1,i}^* := \frac{\beta_i^*}{s_i^3 (\boldsymbol{\alpha}^T s_i \overline{\mathbf{w}_i^*})^2}$ and $\hat{\beta}_{1,i} := \frac{\hat{\beta}_i}{s_i^3 (\boldsymbol{\alpha}^T \mathbf{V}\hat{u}_i)^2}$, and then (B.97) and (B.98) become

$$\hat{\beta}_{1,i} = m_{3,i}(\hat{a}_i), \quad \beta_{1,i}^* = m_{3,i}(\|\mathbf{w}_i^*\|). \quad (\text{B.99})$$

By Assumption 2 and (B.97), there exists a constant $\delta' > 0$ such that the inverse function $g(\cdot)$ of $m_{3,1}(\cdot)$ has upper-bounded derivative in the interval $(\beta_{1,i}^* - \delta', \beta_{1,i}^* + \delta')$, i.e., $|g'(x)| < \Gamma$ for a constant Γ . By employing the result in [7], if the sample size $n \geq d \text{poly}(K, \kappa, t, \log d)$, then \widehat{Q}_1 and Q_1 , $\mathbf{V}\hat{u}_i$ and $s_i \overline{\mathbf{w}_i^*}$ can be arbitrarily close so that $|\beta_{1,i}^* - \hat{\beta}_{1,i}| < \min\{\delta', \frac{r}{\sqrt{K}\Gamma}\}$.

Thus, by (B.99) and the mean value theorem, we obtain

$$|\hat{a}_i - \|\mathbf{w}_i^*\|| = |g'(\xi)| |\beta_{1,i}^* - \hat{\beta}_{1,i}|, \quad (\text{B.100})$$

where ξ is between $\beta_{1,i}^*$ and $\hat{\beta}_{1,i}$, and hence $|g'(\xi)| < \Gamma$. Therefore, $|\hat{a}_i - \|\mathbf{w}_i^*\|| \leq \frac{r}{\sqrt{K}}$, which is the desired result.

Appendix C: Proofs for Chapter 4

C.1 Preliminary

For convenience, we introduce some useful results, for which the proofs of our main result will rely on.

Lemma 16. *Let $\mathbf{z} \in \mathbb{R}^d$ be a standard Gaussian random vector, then for any fixed unit vector $\mathbf{u} \in \mathbb{R}^d$, i.e., $\|\mathbf{u}\|_2 = 1$ and a non-zero fixed vector $\mathbf{w} \in \mathbb{R}^d$, we have the following identities,*

$$\mathbb{E} \left[\mathbf{u}^\top \mathbf{z} \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] = \frac{1}{\sqrt{2\pi}} \frac{\mathbf{u}^\top \mathbf{w}}{\|\mathbf{w}\|_2}, \quad (\text{C.1})$$

$$\mathbb{E} \left[(\mathbf{u}^\top \mathbf{z})^2 \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] = \frac{1}{2}. \quad (\text{C.2})$$

Proof. The basic idea is to decompose \mathbf{u} into an orthogonal pair, i.e., $\mathbf{u} = \mathbf{u}^\perp + \mathbf{u}^\parallel$, where $\mathbf{u}^\parallel = \frac{\mathbf{w}\mathbf{w}^\top}{\|\mathbf{w}\|_2^2} \mathbf{u}$ is the projection of \mathbf{u} onto \mathbf{w} , and $\mathbf{u}^\perp = \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^\top}{\|\mathbf{w}\|_2^2} \right) \mathbf{u}$ is the projection of \mathbf{u} onto the corresponding complementary subspace. Then for (C.1) we have

$$\begin{aligned} \mathbb{E} \left[\mathbf{u}^\top \mathbf{z} \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] &= \mathbb{E} \left[\mathbf{z}^\top (\mathbf{u}^\perp + \mathbf{u}^\parallel) \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] \\ &= \mathbb{E} \left[\mathbf{z}^\top \mathbf{u}^\parallel \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] \\ &= \frac{\mathbf{w}^\top \mathbf{u}}{\|\mathbf{w}\|_2^2} \cdot \mathbb{E} \left[\mathbf{z}^\top \mathbf{w} \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right], \end{aligned}$$

where for the second equality we have used the fact that $\mathbf{z}^\top \mathbf{u}^\perp$ is a zero mean Gaussian random variable and it is independent of $\mathbf{z}^\top \mathbf{w}$. Let $x = \mathbf{w}^\top \mathbf{z} \sim \mathcal{N}(0, \sigma^2)$, where $\sigma^2 = \|\mathbf{w}\|_2^2$, then the result follows by evaluating the expectation:

$$\mathbb{E} \left[\mathbf{z}^\top \mathbf{w} \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] = \int_0^\infty x \cdot \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx = \frac{\sigma}{\sqrt{2\pi}} = \frac{\|\mathbf{w}\|_2}{\sqrt{2\pi}}.$$

Similarly, we can calculate (C.2) as

$$\begin{aligned} \mathbb{E} \left[(\mathbf{u}^\top \mathbf{z})^2 \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] &= \mathbb{E} \left[\left((\mathbf{z}^\top \mathbf{u}^\parallel)^2 + (\mathbf{z}^\top \mathbf{u}^\perp)^2 \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] \\ &= \mathbb{E} \left[(\mathbf{z}^\top \mathbf{u}^\parallel)^2 \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] + \mathbb{E} \left[(\mathbf{z}^\top \mathbf{u}^\perp)^2 \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] \end{aligned} \quad (\text{C.3})$$

$$\begin{aligned} &= \left(\frac{\mathbf{w}^\top \mathbf{u}}{\|\mathbf{w}\|_2^2} \right)^2 \cdot \mathbb{E} \left[(\mathbf{z}^\top \mathbf{w})^2 \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] \\ &\quad + \mathbb{E} \left[(\mathbf{z}^\top \mathbf{u}^\perp)^2 \right] \cdot \mathbb{E} \left[1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right], \end{aligned} \quad (\text{C.4})$$

and calculate the three expectations as follows:

$$\begin{aligned} \mathbb{E} \left[(\mathbf{z}^\top \mathbf{w})^2 \cdot 1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] &= \int_0^\infty x^2 \cdot \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \cdot \sigma^2 \cdot \int_0^\infty \exp\left(-\frac{x^2}{2\sigma^2}\right) dx \\ &= \frac{\sigma^2}{2} \\ &= \frac{\|\mathbf{w}\|_2^2}{2}, \end{aligned}$$

and

$$\begin{aligned} \mathbb{E} \left[1_{\{\mathbf{w}^\top \mathbf{z} > 0\}} \right] &= \mathbb{P}(\mathbf{w}^\top \mathbf{z} > 0) = \frac{1}{2}, \\ \mathbb{E} \left[(\mathbf{z}^\top \mathbf{u}^\perp)^2 \right] &= \|\mathbf{u}^\perp\|_2^2 = 1 - \frac{(\mathbf{u}^\top \mathbf{w})^2}{\|\mathbf{w}\|_2^2}. \end{aligned}$$

Hence, the result follows by plugging the above results to (C.4). \square

Lemma 17 (Lemma 3.1 in [86]). Let $\mathbf{x}_i \in \mathbb{R}^m$ be i.i.d. Gaussian random vectors, then for any $0 < \epsilon < 1$, if $n > c_0 m \epsilon^{-2}$,

$$(1 - \epsilon) \|\mathbf{h}\|_2^2 \leq \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{h})^2 \leq (1 + \epsilon) \|\mathbf{h}\|_2^2 \quad (\text{C.5})$$

holds for all non-zero vectors $\mathbf{h} \in \mathbb{R}^m$ with probability at least $1 - 2 \exp(-c_1 \epsilon^2 n)$, and c_0, c_1 are some constants.

Lemma 18. Denote $\mathbf{h} = \mathbf{w} - \mathbf{w}^*$, and let $\rho = \frac{\mathbf{w}^\top \mathbf{w}^*}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2}$. If $\|\mathbf{h}\|_2 = \|\mathbf{w} - \mathbf{w}^*\|_2 = \nu \cdot \|\mathbf{w}^*\|_2$ where $0 \leq \nu < 1$, then we have

$$\rho \geq 1 - \nu^2 \quad (\text{C.6})$$

holds.

Proof. We first show that the assumption $\nu < 1$ implies $\rho > 0$ since

$$\|\mathbf{h}\|_2^2 = \|\mathbf{w}^*\|_2^2 + \|\mathbf{w}\|_2^2 - 2\rho \|\mathbf{w}^*\|_2 \|\mathbf{w}\|_2 < \|\mathbf{w}^*\|_2^2, \quad (\text{C.7})$$

gives us

$$\rho > \frac{\|\mathbf{w}\|_2}{2\|\mathbf{w}^*\|_2} > 0. \quad (\text{C.8})$$

We next show a tighter lower bound of ρ . Let α be the angle between \mathbf{h} and \mathbf{w}^* , i.e.,

$\cos(\alpha) = \frac{\mathbf{h}^\top \mathbf{w}^*}{\|\mathbf{h}\|_2 \|\mathbf{w}^*\|_2}$. We can calculate ρ as

$$\rho = \frac{\mathbf{w}^\top \mathbf{w}^*}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2} = \frac{\mathbf{h}^\top \mathbf{w}^* + \|\mathbf{w}^*\|_2^2}{\|\mathbf{w}^* + \mathbf{h}\|_2 \|\mathbf{w}^*\|_2} = \frac{\cos(\alpha) \|\mathbf{h}\|_2 + \|\mathbf{w}^*\|_2}{\|\mathbf{h} + \mathbf{w}^*\|_2},$$

square the two sides will give us,

$$\rho^2 = \frac{\cos(\alpha)^2 \cdot \|\mathbf{h}\|_2^2 + \|\mathbf{w}^*\|_2^2 + 2\cos(\alpha) \|\mathbf{h}\|_2 \|\mathbf{w}^*\|_2}{\|\mathbf{w}^*\|_2^2 + \|\mathbf{h}\|_2^2 + 2\cos(\alpha) \|\mathbf{w}^*\|_2 \|\mathbf{h}\|_2} = 1 - \frac{(1 - \cos^2(\alpha)) \cdot \nu^2}{1 + \nu^2 + 2\nu \cos(\alpha)}. \quad (\text{C.9})$$

Moreover, when $\cos(\alpha)^2 \neq 1$,

$$\frac{(1 - \cos^2(\alpha)) \cdot \nu^2}{1 + \nu^2 + 2\nu\cos(\alpha)} = \frac{\nu^2}{1 + \frac{(\nu + \cos(\alpha))^2}{1 - \cos^2(\alpha)}} \leq \nu^2,$$

where the equality holds when $\cos(\alpha) = -\nu$, and when $\cos(\alpha)^2 = 1$, $\rho^2 = 1$. Together we can conclude that

$$1 \geq \rho \geq \rho^2 \geq 1 - \nu^2. \quad (\text{C.10})$$

□

C.2 Proof of GD on the Population Risk

C.2.1 Proof of Theorem 6

Proof. According to the definition of strong convexity, we can calculate the Hessian of the population risk function and verify it is positive definite in certain region. We next derive the Hessian and check its spectrum. First of all, according to [71], we can write the population risk (up to additive factors in \mathbf{w}^*) and its gradient as

$$\begin{aligned} L(\mathbf{w}) &= \frac{1}{K^2} \left[\left(\frac{K^2 - K}{2\pi} + \frac{K}{2} \right) \|\mathbf{w}\|_2^2 - \frac{K}{\pi} \|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2 (\sin(\theta) + (\pi - \theta) \cos(\theta)) \right. \\ &\quad \left. - \frac{K^2 - K}{\pi} \|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2 \right], \\ \nabla L(\mathbf{w}) &= \left(\frac{1}{K} + \frac{1 - \frac{1}{K}}{\pi} \right) \cdot \mathbf{w} - \frac{\|\mathbf{w}^*\|_2}{\pi K} \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \cdot \sin(\theta) - \frac{1 - \frac{1}{K}}{\pi} \cdot \|\mathbf{w}^*\|_2 \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \\ &\quad - \frac{\pi - \theta}{K\pi} \cdot \mathbf{w}^*. \end{aligned}$$

When $\theta \neq 0$ or π and $\mathbf{w} \neq \mathbf{0}$, we can calculate the Hessian as

$$\begin{aligned}
& \nabla^2 L(\mathbf{w}) \\
&= \left(\frac{1}{K} + \frac{1 - \frac{1}{K}}{\pi} \right) \cdot \mathbf{I} - \frac{\|\mathbf{w}^*\|_2}{\pi K} \cdot \left[\frac{\mathbf{w}}{\|\mathbf{w}\|_2} \cdot \cos\theta \cdot \frac{\partial\theta}{\partial\mathbf{w}} + \sin\theta \left(\frac{1}{\|\mathbf{w}\|_2} \cdot \mathbf{I} - \frac{1}{\|\mathbf{w}\|_2^3} \cdot \mathbf{w}\mathbf{w}^\top \right) \right] \\
&\quad - \frac{1 - \frac{1}{K}}{\pi} \cdot \|\mathbf{w}^*\|_2 \left(\frac{1}{\|\mathbf{w}\|_2} \cdot \mathbf{I} - \frac{1}{\|\mathbf{w}\|_2^3} \cdot \mathbf{w}\mathbf{w}^\top \right) + \frac{1}{K\pi} \cdot \mathbf{w}^* \cdot \frac{\partial\theta}{\partial\mathbf{w}} \\
&= \underbrace{\left(\frac{1}{K} + \frac{1 - \frac{1}{K}}{\pi} - \frac{\sin(\theta)}{\pi K} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} - \frac{1 - \frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \right) \cdot \mathbf{I}}_{\mathbf{H}_1} \\
&\quad + \underbrace{\frac{1}{\pi K} \cdot \frac{\cos(\theta)}{\sin(\theta)} \cdot \frac{1}{\|\mathbf{w}\|_2^2} \cdot (\mathbf{w}^*\mathbf{w}^\top + \mathbf{w}\mathbf{w}^{*\top})}_{\mathbf{H}_2} \\
&\quad + \underbrace{\left(\frac{1 - \frac{1}{K}}{\pi} - \frac{1}{\pi K} \cdot \frac{\cos(2\theta)}{\sin(\theta)} \right) \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2^3} \cdot \mathbf{w}\mathbf{w}^\top}_{\mathbf{H}_3} \\
&\quad - \underbrace{\frac{1}{K\pi} \frac{1}{\sin(\theta)} \cdot \frac{1}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2} \mathbf{w}^*\mathbf{w}^{*\top}}_{\mathbf{H}_4}. \tag{C.11}
\end{aligned}$$

It is easy to show that the first term \mathbf{H}_1 is positively definite under some mild conditions. The difficulty actually lies in lower bounding the spectrum of the last three terms. We then lower bound the smallest eigenvalue of $\mathbf{H}_2 + \mathbf{H}_3 + \mathbf{H}_4$ as

following,

$$\begin{aligned}
& \min_{\|\mathbf{u}\|_2=1} \mathbf{u}^\top (\mathbf{H}_2 + \mathbf{H}_3 + \mathbf{H}_4) \mathbf{u} \\
&= \min_{\|\mathbf{u}\|_2=1} \frac{2}{\pi K} \cdot \frac{\cos(\theta)}{\sin(\theta)} \cdot \frac{1}{\|\mathbf{w}\|_2^2} \cdot \mathbf{u}^\top \mathbf{w} \cdot \mathbf{u}^\top \mathbf{w}^* - \frac{1}{K\pi} \frac{1}{\sin(\theta)} \cdot \frac{(\mathbf{u}^\top \mathbf{w}^*)^2}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2} \\
&\quad + \left(\frac{1 - \frac{1}{K}}{\pi} - \frac{1}{\pi K} \cdot \frac{\cos(2\theta)}{\sin(\theta)} \right) \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2^3} \cdot (\mathbf{u}^\top \mathbf{w})^2 \\
&= \min_{\|\mathbf{u}\|_2=1} \frac{2}{\pi K} \cdot \frac{\cos(\theta)}{\sin(\theta)} \cdot \frac{1}{\|\mathbf{w}\|_2^2} \cdot \mathbf{u}^\top \mathbf{w} \cdot \mathbf{u}^\top \mathbf{w}^* - \frac{1}{\pi K} \cdot \frac{\cos(2\theta)}{\sin(\theta)} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2^3} \cdot (\mathbf{u}^\top \mathbf{w})^2 \\
&\quad - \frac{1}{K\pi} \frac{1}{\sin(\theta)} \cdot \frac{(\mathbf{u}^\top \mathbf{w}^*)^2}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2} + \frac{1 - \frac{1}{K}}{\pi} \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2^3} \cdot (\mathbf{u}^\top \mathbf{w})^2 \\
&= \min_{\|\mathbf{u}\|_2=1} \frac{1}{\pi K} \cdot \frac{1}{\sin(\theta)} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \left[2\cos(\theta) \mathbf{u}^\top \bar{\mathbf{w}} \cdot \mathbf{u}^\top \bar{\mathbf{w}}^* - \cos(\theta)^2 (\mathbf{u}^\top \bar{\mathbf{w}})^2 \right. \\
&\quad \left. + \sin(\theta)^2 (\mathbf{u}^\top \bar{\mathbf{w}})^2 - (\mathbf{u}^\top \bar{\mathbf{w}}^*)^2 \right] + \frac{1 - \frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} (\mathbf{u}^\top \bar{\mathbf{w}})^2 \\
&= \min_{\|\mathbf{u}\|_2=1} \frac{1}{\pi K} \cdot \frac{1}{\sin(\theta)} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \left[\sin(\theta)^2 (\mathbf{u}^\top \bar{\mathbf{w}})^2 - (\mathbf{u}^\top (\mathbf{I} - \bar{\mathbf{w}}\bar{\mathbf{w}}^\top) \bar{\mathbf{w}}^*)^2 \right. \\
&\quad \left. + (K - 1) \sin(\theta) (\mathbf{u}^\top \bar{\mathbf{w}})^2 \right] \tag{C.12} \\
&\geq \frac{1}{\pi K} \cdot \frac{1}{\sin(\theta)} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \cdot (-\sin(\theta)^2) \\
&= -\frac{\sin(\theta)}{\pi K} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2}
\end{aligned}$$

where we use $\bar{\mathbf{w}}$ to denote $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, and note that $\bar{\mathbf{w}} \perp (\mathbf{I} - \bar{\mathbf{w}}\bar{\mathbf{w}}^\top) \bar{\mathbf{w}}^*$, hence setting \mathbf{u} in the direction of $(\mathbf{I} - \bar{\mathbf{w}}\bar{\mathbf{w}}^\top) \bar{\mathbf{w}}^*$ will achieve the minimum of (C.12). Moreover,

$$\max_{\|\mathbf{u}\|_2=1} \mathbf{u}^\top (\mathbf{H}_2 + \mathbf{H}_3 + \mathbf{H}_4) \mathbf{u} \leq \frac{1}{\pi} \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \tag{C.13}$$

holds by setting \mathbf{u} in the direction of $\bar{\mathbf{w}}$. Hence if the condition $\|\mathbf{w}\|_2 \geq \frac{K+1}{K-1+\frac{3}{4}\pi} \|\mathbf{w}^*\|_2$ holds, we have the eigenvalue of the Hessian is lower bounded by some positive value,

$$\begin{aligned}
& \min_{\|\mathbf{u}\|_2=1} \mathbf{u}^\top \mathbf{H}_1 \mathbf{u} + \mathbf{u}^\top (\mathbf{H}_2 + \mathbf{H}_3 + \mathbf{H}_4) \mathbf{u} \\
& \geq \frac{1}{K} + \frac{1 - \frac{1}{K}}{\pi} - \frac{\sin(\theta)}{\pi K} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} - \frac{1 - \frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} - \frac{\sin(\theta)}{\pi K} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \\
& \geq \frac{1}{K} + \frac{1 - \frac{1}{K}}{\pi} - \left(\frac{2}{\pi K} + \frac{K-1}{K\pi} \right) \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \\
& \geq \frac{1}{4K},
\end{aligned} \tag{C.14}$$

and we can also obtain the upper bound as

$$\begin{aligned}
& \max_{\|\mathbf{u}\|_2=1} \mathbf{u}^\top (\mathbf{H}_1) \mathbf{u} + \mathbf{u}^\top (\mathbf{H}_2 + \mathbf{H}_3 + \mathbf{H}_4) \mathbf{u} \\
& \leq \frac{1}{K} + \frac{1 - \frac{1}{K}}{\pi} - \frac{\sin(\theta)}{\pi K} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} - \frac{1 - \frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} + \frac{1}{\pi} \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \\
& \leq \frac{K-1+\pi}{\pi K} + \frac{1}{\pi K} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2} \\
& \leq 3.
\end{aligned} \tag{C.15}$$

As a summary, when $\theta \neq 0$ or π and $\mathbf{w} \neq 0$ we have shown that

$$\frac{1}{4K} \cdot \mathbf{I} \preceq \nabla^2 L(\mathbf{w}) \preceq 3 \cdot \mathbf{I} \tag{C.16}$$

holds when $\|\mathbf{w}\|_2 \geq \frac{K+1}{K-1+\frac{3}{4}\pi} \|\mathbf{w}^*\|_2$.

Next we will consider the case $\theta = 0$ and calculate the Hessian i.e., the derivative of $\nabla L(\mathbf{w})$ at $\lambda \cdot \mathbf{w}^*$ where $\lambda \in \mathbb{R}^+$. We decompose $\nabla L(\mathbf{w})$ as two terms

$$\nabla L(\mathbf{w}) = g_1(\mathbf{w}) + g_2(\mathbf{w}), \tag{C.17}$$

where

$$g_1(\mathbf{w}) = \frac{\|\mathbf{w}^*\|_2}{K\pi} \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \cdot \sin(\theta) - \frac{\theta}{K\pi} \cdot \mathbf{w}^*, \tag{C.18}$$

$$g_2(\mathbf{w}) = \left(\frac{1}{K} + \frac{1 - \frac{1}{K}}{\pi} \right) \cdot \mathbf{w} - \frac{1 - \frac{1}{K}}{\pi} \cdot \|\mathbf{w}^*\|_2 \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|_2}. \tag{C.19}$$

For $g_1(\mathbf{w})$, according to the definition [87, Definition 9.11], if there exist $\mathbf{A} \in \mathbb{R}^{m \times m}$ such that

$$\lim_{\epsilon \rightarrow 0} \frac{\|g_1(\lambda\mathbf{w}^* + \epsilon\mathbf{u}) - \mathbf{A} \cdot \epsilon \cdot \mathbf{u}\|_2}{\|\epsilon\mathbf{u}\|_2} = 0, \quad (\text{C.20})$$

holds for all $\|\mathbf{u}\|_2 = 1$, then we say that $g_1(\mathbf{w})$ is differentiable at $\lambda\mathbf{w}^*$ and we write $\nabla g_1(\lambda\mathbf{w}^*) = \mathbf{A}$. Plugging $g_1(\mathbf{w})$ into the left-hand side of (C.20) we obtain

$$\begin{aligned} & \lim_{\epsilon \rightarrow 0} \left\| \frac{\frac{\|\mathbf{w}^*\|_2}{K\pi} \cdot \frac{\lambda\mathbf{w}^* + \epsilon\mathbf{u}}{\|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2} \cdot \sin(\theta_\epsilon) - \frac{\theta_\epsilon}{\pi K} \cdot \mathbf{w}^*}{\epsilon} - \mathbf{A}\mathbf{u} \right\|_2 \\ &= \lim_{\epsilon \rightarrow 0} \left\| \frac{\|\mathbf{w}^*\|_2}{K\pi} \cdot \frac{\lambda\mathbf{w}^*}{\|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2} \cdot \frac{\sin(\theta_\epsilon)}{\epsilon} + \frac{\|\mathbf{w}^*\|_2}{K\pi} \cdot \frac{\mathbf{u}}{\|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2} \cdot \sin(\theta_\epsilon) \right. \\ & \quad \left. - \frac{\theta_\epsilon}{\epsilon} \cdot \frac{\mathbf{w}^*}{\pi K} - \mathbf{A} \cdot \mathbf{u} \right\|_2 \end{aligned}$$

where θ_ϵ is the angle between $\lambda\mathbf{w}^* + \epsilon\mathbf{u}$ and $\lambda\mathbf{w}^*$, and we have

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \frac{\sin(\theta_\epsilon)}{\epsilon} &= \lim_{\epsilon \rightarrow 0} \frac{\sqrt{1 - \left(\frac{\lambda\mathbf{w}^{*\top}(\lambda\mathbf{w}^* + \epsilon\mathbf{u})}{\|\lambda\mathbf{w}^*\|_2 \|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2} \right)^2}}{\epsilon} = \lim_{\epsilon \rightarrow 0} \sqrt{\frac{\|\lambda\mathbf{w}^*\|_2^2 - (\lambda\mathbf{w}^{*\top} \mathbf{u})^2}{\|\lambda\mathbf{w}^*\|_2^2 \|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2^2}} \\ &= \sqrt{\frac{\|\lambda\mathbf{w}^*\|_2^2 - (\lambda\mathbf{w}^{*\top} \mathbf{u})^2}{\|\lambda\mathbf{w}^*\|_2^2 \|\lambda\mathbf{w}^*\|_2^2}} \quad (\text{C.21}) \end{aligned}$$

exists. Further since,

$$\lim_{\epsilon \rightarrow 0} \frac{\theta_\epsilon}{\epsilon} = \lim_{\epsilon \rightarrow 0} \left(\frac{\theta_\epsilon}{\sin(\theta_\epsilon)} \cdot \frac{\sin(\theta_\epsilon)}{\epsilon} \right) \quad (\text{C.22})$$

and

$$\lim_{\epsilon \rightarrow 0} \frac{\theta_\epsilon}{\sin(\theta_\epsilon)} = \lim_{\epsilon \rightarrow 0} \frac{1}{\cos(\theta_\epsilon)} = \lim_{\epsilon \rightarrow 0} \frac{\|\lambda\mathbf{w}^*\|_2 \|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2}{\lambda\mathbf{w}^{*\top}(\lambda\mathbf{w}^* + \epsilon\mathbf{u})} = 1, \quad \lim_{\epsilon \rightarrow 0} \frac{\theta_\epsilon}{\epsilon} = \lim_{\epsilon \rightarrow 0} \frac{\sin(\theta_\epsilon)}{\epsilon}, \quad (\text{C.23})$$

both exist, hence

$$\lim_{\epsilon \rightarrow 0} \frac{\|\mathbf{w}^*\|_2}{K\pi} \cdot \frac{\lambda\mathbf{w}^*}{\|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2} \cdot \frac{\sin(\theta_\epsilon)}{\epsilon} - \frac{\theta_\epsilon}{\epsilon} \cdot \frac{\mathbf{w}^*}{\pi K} = 0. \quad (\text{C.24})$$

Together with

$$\lim_{\epsilon \rightarrow 0} \frac{\|\mathbf{w}^*\|_2}{K\pi} \cdot \frac{\mathbf{u}}{\|\lambda\mathbf{w}^* + \epsilon\mathbf{u}\|_2} \cdot \sin(\theta_\epsilon) = 0, \quad (\text{C.25})$$

we will have (C.20) hold when $\mathbf{A} = \mathbf{0}$, i.e.,

$$\nabla g_1(\lambda\mathbf{w}^*) = \mathbf{0}. \quad (\text{C.26})$$

Next for $g_2(\mathbf{w}) = \left(\frac{1}{K} + \frac{1-\frac{1}{K}}{\pi}\right) \cdot \mathbf{w} - \frac{1-\frac{1}{K}}{\pi} \cdot \|\mathbf{w}^*\|_2 \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, the derivative can be calculated as

$$\nabla g_2(\mathbf{w}) = \left(\frac{1}{K} + \frac{1-\frac{1}{K}}{\pi} - \frac{1-\frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2}\right) \cdot \mathbf{I} + \frac{1-\frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2^3} \cdot \mathbf{w}\mathbf{w}^\top, \quad (\text{C.27})$$

and since $\nabla L(\mathbf{w}) = g_1(\mathbf{w}) + g_2(\mathbf{w})$, we can conclude that when $\theta = 0$,

$$\begin{aligned} \nabla^2 L(\mathbf{w}) &= \left(\frac{1}{K} + \frac{1-\frac{1}{K}}{\pi} - \frac{1-\frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2}\right) \cdot \mathbf{I} + \frac{1-\frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2^3} \cdot \mathbf{w}\mathbf{w}^\top \\ &\succeq \left(\frac{1}{K} + \frac{1-\frac{1}{K}}{\pi} - \frac{1-\frac{1}{K}}{\pi} \cdot \frac{\|\mathbf{w}^*\|_2}{\|\mathbf{w}\|_2}\right) \cdot \mathbf{I}, \end{aligned} \quad (\text{C.28})$$

and

$$\nabla^2 L(\mathbf{w}) \preceq \left(\frac{1}{K} + \frac{1-\frac{1}{K}}{\pi}\right) \cdot \mathbf{I}. \quad (\text{C.29})$$

Hence when $\theta = 0$, the result still holds.

□

C.3 Proof of GD on the Empirical Risk

C.3.1 Proof of Lemma 1

Proof. Recall that the empirical risk function (4.3) is

$$L_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}^\top \mathbf{x}_i^{(j)}) - y_i \right)^2, \quad (\text{C.30})$$

and we can calculate its gradient as

$$\begin{aligned}\nabla L_n(\mathbf{w}) &= 2 \cdot \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}^\top \mathbf{x}_i^{(j)}) - y_i \right) \cdot \left(\frac{1}{K} \sum_{j=1}^K \phi'(\mathbf{w}^\top \mathbf{x}_i^{(j)}) \mathbf{x}_i^{(j)} \right) \\ &= 2 \cdot \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \left(\phi(\mathbf{w}^\top \mathbf{x}_i^{(j)}) - \phi(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}) \right) \phi'(\mathbf{w}^\top \mathbf{x}_i^{(l)}) \mathbf{x}_i^{(l)},\end{aligned}\tag{C.31}$$

where ϕ' denotes the derivative of ϕ , and following the convention we define it as

$$\phi'(z) = 1_{\{z>0\}} = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}.\tag{C.32}$$

Further denote $\mathbf{h} = \mathbf{w} - \mathbf{w}^*$, then the left-hand side of (4.15) can be written as

$$\begin{aligned}\langle \nabla L_n(\mathbf{w}), \mathbf{h} \rangle &= 2 \cdot \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \left[\mathbf{h}^\top \mathbf{x}_i^{(j)} \cdot \mathbf{h}^\top \mathbf{x}_i^{(l)} \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right. \\ &\quad \left. + \mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} \cdot \mathbf{h}^\top \mathbf{x}_i^{(l)} \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0\}} - 1_{\{\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} > 0\}} \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right].\end{aligned}\tag{C.33}$$

One pivotal observation of (C.33) is that when \mathbf{w} is close to \mathbf{w}^* , the first term can be lower bounded by the distance $\|\mathbf{h}\|_2$, and the second term is very small since $\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}$ and $\mathbf{w}^\top \mathbf{x}_i^{(j)}$ will have the same sign with high probability. We summarize the observation in the following two lemmas.

Lemma 19. *When the sample complexity satisfies $n \geq c \cdot mK^2 \cdot \log(n)$ for some sufficiently large constant c , with probability at least $1 - d^{-10}$, we have*

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{h}^\top \mathbf{x}_i^{(j)} \cdot \mathbf{h}^\top \mathbf{x}_i^{(l)} \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \geq \frac{1}{4K} \cdot \|\mathbf{h}\|_2^2\tag{C.34}$$

holds for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, \frac{1}{K} \|\mathbf{w}^*\|_2)$.

Lemma 20. For some sufficiently large constant c_1, c_2 , if $n \geq c_1 \cdot mK^2 \log(K)$ and $\|\mathbf{h}\|_2 \leq \frac{1}{K} \|\mathbf{w}^*\|_2$ hold, then with probability at least $1 - c_2 \exp(-K^2 n)$,

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} \cdot \mathbf{h}^\top \mathbf{x}_i^{(l)} \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0\}} - 1_{\{\mathbf{w}^{*\top} \mathbf{x}_i^{(j)} > 0\}} \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \\ & \leq \frac{1.2}{K^{\frac{3}{2}}} \cdot \|\mathbf{h}\|_2^2 \end{aligned} \quad (\text{C.35})$$

holds for all $\mathbf{w} \in \mathbb{B} \left(\mathbf{w}^*, \frac{1}{K^{\frac{3}{2}}} \|\mathbf{w}^*\|_2 \right)$.

The proof of Lemma 19 and Lemma 20 is provided in Section C.3.2. Together with the above two lemmas we can conclude that as long as $n \geq c \cdot mK^2 \cdot \log(n)$ for some large constant c and $\|\mathbf{h}\|_2 \leq \frac{1}{K^{\frac{3}{2}}} \|\mathbf{w}^*\|_2$,

$$\langle \nabla L_n(\mathbf{w}), \mathbf{h} \rangle \geq \frac{\sqrt{K} - 5.8}{4K^{\frac{3}{2}}} \cdot \|\mathbf{h}\|_2^2 \quad (\text{C.36})$$

holds with probability at least $1 - \frac{1}{d^{10}}$ for all \mathbf{w} . For the right-hand side of RC (4.15) we successfully upper bound $\|\nabla L_n(\mathbf{w})\|_2^2$, and the result is summarized as follows.

Lemma 21. If $n \geq c_1 m \epsilon^{-2}$, then with probability at least $1 - 2 \exp(-c_2 \epsilon^2 n) - 2 \exp(-n/2)$ we have

$$\|\nabla L_n(\mathbf{w})\|_2^2 \leq 18(1 + \epsilon) \|\mathbf{h}\|_2^2 \quad (\text{C.37})$$

holds for all \mathbf{w} .

Setting $\epsilon = \frac{1}{K}$ in lemma 21 we have the following inequality

$$\|\nabla L_n(\mathbf{w})\|_2^2 \leq 18 \left(1 + \frac{1}{K} \right) \|\mathbf{h}\|_2^2 \quad (\text{C.38})$$

holds with probability at least $1 - c \cdot \exp(-mK^2)$ for some large constant c as long as $n \geq c_4 \cdot mK^2$. Hence, further take $\lambda = \frac{1}{72K^2}$ and $\mu = \frac{1}{4K^2}$, then we have

$$\langle \nabla L_n(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \lambda \|\nabla L_n(\mathbf{w})\|_2^2 + \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2 \quad (\text{C.39})$$

holds with probability at least $1 - \frac{1}{d^{10}}$ as long as the sample complexity satisfies $n \geq c \cdot mK^2 \cdot \log(n)$ for some large constant c and $\|\mathbf{h}\|_2 \leq \frac{1}{K^{\frac{3}{2}}}\|\mathbf{w}^*\|_2$. \square

C.3.2 Proof of Auxiliary Lemmas

Proof of Lemma 19

Proof. Lower bounding (C.34) can be boiled down to lower bound the smallest eigenvalue of the symmetric matrix $\mathbf{A}_n(\mathbf{w}) \in \mathbb{R}^{m \times m}$, where

$$\mathbf{A}_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{x}_i^{(j)} \cdot \mathbf{x}_i^{(l)\top} \cdot \mathbf{1}_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}}. \quad (\text{C.40})$$

Due to the randomness of $\mathbf{A}_n(\mathbf{w})$, it's hard to lower bound $\mathbf{A}_n(\mathbf{w})$ directly, however it is easy to calculate and lower bound $\mathbb{E}[\mathbf{A}_n(\mathbf{w})]$. Hence, we first lower bound $\mathbb{E}[\mathbf{A}_n(\mathbf{w})]$, and then show that $\mathbf{A}_n(\mathbf{w})$ concentrates around its expectation.

With the two identities from Lemma 16, we can lower bound $\mathbb{E}[\mathbf{A}_n(\mathbf{w})]$ as

$$\begin{aligned} & \min_{\|\mathbf{u}\|_2=1} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{u}^\top \mathbf{x}_i^{(j)} \cdot \mathbf{x}_i^{(l)\top} \mathbf{u} \cdot \mathbf{1}_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right] \cdot \mathbf{I} \\ &= \min_{\|\mathbf{u}\|_2=1} \frac{1}{K^2} \left(\sum_{j \neq l} \mathbb{E} \left[\mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{1}_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0\}} \right] \cdot \mathbb{E} \left[\mathbf{u}^\top \mathbf{x}^{(l)} \cdot \mathbf{1}_{\{\mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} \right] \right. \\ & \quad \left. + \sum_{j=1}^K \mathbb{E} \left[(\mathbf{u}^\top \mathbf{x}^{(j)})^2 \cdot \mathbf{1}_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0\}} \right] \right) \cdot \mathbf{I} \\ &= \min_{\|\mathbf{u}\|_2=1} \frac{1}{K^2} \left((K^2 - K) \cdot \frac{1}{2\pi} \frac{(\mathbf{u}^\top \mathbf{w})^2}{\|\mathbf{w}\|_2^2} + \frac{K}{2} \right) \cdot \mathbf{I} \\ &\succeq \frac{1}{2K} \cdot \mathbf{I} \end{aligned} \quad (\text{C.41})$$

where in the first equality we have used the fact that $\mathbf{x}^{(j)}$ is independent of $\mathbf{x}^{(l)}$ when $j \neq l$, and for the last inequality, the equality holds when $\langle \mathbf{u}, \mathbf{w} \rangle = 0$. In a summary, we have shown that

$$\mathbb{E}[\mathbf{A}_n(\mathbf{w})] \succeq \frac{1}{2K} \cdot \mathbf{I} \quad (\text{C.42})$$

holds. Next we are going to apply a covering argument to show that with high probability the perturbation $\|\mathbf{A}_n(\mathbf{w}) - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\|$ is small for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, \frac{1}{K}\|\mathbf{w}^*\|_2)$. Notice that only the sign of $\mathbf{w}^\top \mathbf{x}$ matters for $\mathbf{A}_n(\mathbf{w})$, hence considering $\mathbf{w} \in \mathbb{S}^{m-1}$ will be enough. Denote \mathcal{N}_ϵ as the ϵ -net of the sphere \mathbb{S}^{m-1} , i.e. for any $\mathbf{w} \in \mathbb{S}^{m-1}$, there exist a $\mathbf{w}_\epsilon \in \mathcal{N}_\epsilon$ corresponding to \mathbf{w} such that $\|\mathbf{w} - \mathbf{w}_\epsilon\|_2 \leq \epsilon$. By triangle inequality, we can write

$$\begin{aligned} \|\mathbf{A}_n(\mathbf{w}) - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| &\leq \|\mathbf{A}_n(\mathbf{w}) - \mathbf{A}_n(\mathbf{w}_\epsilon)\| + \|\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]\| \\ &\quad + \|\mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)] - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\|, \end{aligned} \quad (\text{C.43})$$

and hence

$$\begin{aligned} &\mathbb{P}\left(\sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \|\mathbf{A}_n(\mathbf{w}) - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| \geq t\right) \\ &\leq \mathbb{P}\left(\sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \|\mathbf{A}_n(\mathbf{w}) - \mathbf{A}_n(\mathbf{w}_\epsilon)\| \geq \frac{t}{3}\right) \\ &\quad + \mathbb{P}\left(\sup_{\mathbf{w}_\epsilon \in \mathcal{N}_\epsilon} \|\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]\| \geq \frac{t}{3}\right) \\ &\quad + \mathbb{P}\left(\sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \|\mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)] - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| \geq \frac{t}{3}\right). \end{aligned} \quad (\text{C.44})$$

Next we will deal with the above terms step by step.

Firstly, let $V_{\frac{1}{4}}$ be a $\frac{1}{4}$ -net of $\mathbb{B}(0, 1)$ with $\log |V_{\frac{1}{4}}| \leq m \log 12$. We have

$$\|\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]\| \leq 2 \sup_{\mathbf{u} \in V_{\frac{1}{4}}} |\langle \mathbf{u}, (\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]) \mathbf{u} \rangle|, \quad (\text{C.45})$$

and then applying union bound over N_ϵ and $V_{\frac{1}{4}}$ will give us

$$\begin{aligned}
& \mathbb{P} \left(\sup_{\mathbf{w}_\epsilon \in N_\epsilon} \|\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]\| \geq \frac{t}{3} \right) \\
& \leq |N_\epsilon| \cdot |V_{\frac{1}{4}}| \cdot \sup_{\mathbf{w}_\epsilon \in N_\epsilon} \sup_{\mathbf{u} \in V_{\frac{1}{4}}} \mathbb{P} \left(|\langle \mathbf{u}, (\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]) \mathbf{u} \rangle| \geq \frac{t}{6} \right) \\
& \leq \exp \left(m \log \left(1 + \frac{2}{\epsilon} \right) + m \log(12) \right) \\
& \quad \cdot \sup_{\mathbf{w}_\epsilon \in N_\epsilon} \sup_{\mathbf{u} \in V_{\frac{1}{4}}} \mathbb{P} \left(|\langle \mathbf{u}, (\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]) \mathbf{u} \rangle| \geq \frac{t}{6} \right). \tag{C.46}
\end{aligned}$$

Further since the sub-exponential norm of $\mathbf{A}_n(\mathbf{w})$ can be upper bounded as

$$\begin{aligned}
& \left\| \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u} \cdot \mathbf{1}_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} \right\|_{\psi_1} \\
& \leq \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \|\mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u}\|_{\psi_1} \\
& \leq C,
\end{aligned}$$

where C is some constant. Hence, applying the Bernstein inequality will give us

$$\mathbb{P} \left(|\langle \mathbf{u}, (\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]) \mathbf{u} \rangle| \geq \frac{t}{6} \right) \leq 2 \exp(-C \cdot n \cdot \min\{c_1 \cdot t^2, c_2 \cdot t\}), \tag{C.47}$$

together we will obtain

$$\begin{aligned}
& \mathbb{P} \left(\sup_{\mathbf{w}_\epsilon \in N_\epsilon} \|\mathbf{A}_n(\mathbf{w}_\epsilon) - \mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)]\| \geq \frac{t}{3} \right) \\
& \leq \exp \left(m \log \left(12 + \frac{24}{\epsilon} \right) \right) \cdot 2 \exp(-C \cdot n \cdot \min\{c_1 \cdot t^2, c_2 \cdot t\}).
\end{aligned}$$

Secondly, by the definition of spectral norm and apply Lemma 16 will give us

$$\begin{aligned}
\|\mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)] - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| &= \sup_{\|\mathbf{u}\|_2=1} \mathbf{u}^\top (\mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)] - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]) \mathbf{u} \\
&= \sup_{\|\mathbf{u}\|_2=1} \frac{K-1}{2\pi K} \left(\frac{(\mathbf{u}^\top \mathbf{w}_\epsilon)^2}{\|\mathbf{w}_\epsilon\|_2^2} - \frac{(\mathbf{u}^\top \mathbf{w})^2}{\|\mathbf{w}\|_2^2} \right), \tag{C.48}
\end{aligned}$$

recall that $\|\mathbf{w}\|_2 = \|\mathbf{w}_\epsilon\|_2 = 1$ and plugging it back to the above equation,

$$\begin{aligned} \|\mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)] - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| &= \sup_{\|\mathbf{u}\|_2=1} \frac{K-1}{2\pi K} (\mathbf{u}^\top (\mathbf{w}_\epsilon - \mathbf{w}) \cdot \mathbf{u}^\top (\mathbf{w}_\epsilon + \mathbf{w})) \\ &< \frac{K-1}{2\pi K} \|\mathbf{w}_\epsilon - \mathbf{w}\|_2 \cdot \|\mathbf{w}_\epsilon + \mathbf{w}\|_2, \end{aligned} \quad (\text{C.49})$$

where in the last step we apply Cauchy-Schwarz inequality. Hence by the definition of \mathbf{w}_ϵ we have the following holds

$$\|\mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)] - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| < \frac{K-1}{\pi K} \cdot \epsilon. \quad (\text{C.50})$$

Thus if we set $\frac{t}{3} \geq \frac{K-1}{\pi K} \cdot \epsilon$ we will have

$$\mathbb{P}\left(\sup_{\mathbf{w} \in N_\epsilon} \|\mathbb{E}[\mathbf{A}_n(\mathbf{w}_\epsilon)] - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| \geq \frac{t}{3}\right) = 0. \quad (\text{C.51})$$

Finally, by Markov's inequality we will obtain

$$\mathbb{P}\left(\sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \|\mathbf{A}_n(\mathbf{w}) - \mathbf{A}_n(\mathbf{w}_\epsilon)\| \geq \frac{t}{3}\right) \leq \frac{3}{t} \cdot \mathbb{E}\left[\sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \|\mathbf{A}_n(\mathbf{w}) - \mathbf{A}_n(\mathbf{w}_\epsilon)\|\right], \quad (\text{C.52})$$

and we can upper bound the right-hand side as following

$$\begin{aligned} &\mathbb{E}\left[\sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \|\mathbf{A}_n(\mathbf{w}) - \mathbf{A}_n(\mathbf{w}_\epsilon)\|\right] \\ &= \mathbb{E}\left[\sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \left\| \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \right. \right. \\ &\quad \left. \left. \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}}\right) \right\|\right] \\ &\leq \mathbb{E}\left[\left\| \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \right\| \right. \end{aligned} \quad (\text{C.53})$$

$$\begin{aligned} &\left. \cdot \left(\max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} |1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}}| \right) \right] \\ &\leq \sqrt{\mathbb{E}\left[\left\| \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \right\|^2\right]} \end{aligned} \quad (\text{C.54})$$

$$\cdot \sqrt{\mathbb{E}\left[\left(\max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} |1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}}| \right)^2\right]}, \quad (\text{C.55})$$

the last inequality follows from Cauchy-Schwarz inequality. Notice that

$$\begin{aligned}
& \left(\max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \left| 1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}} \right| \right)^2 \\
&= \max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \left| 1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}} \right|^2, \\
&= \max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \left| 1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}} \right|
\end{aligned} \tag{C.56}$$

since $|1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}}|$ is either 0 or 1. Applying triangle inequality we have

$$\begin{aligned}
& \max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \left| 1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}} \right| \\
& \leq \max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \left| 1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0\}} \right| + \left| 1_{\{\mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}} \right|
\end{aligned} \tag{C.57}$$

hold. Furthermore, we have

$$\left| 1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0\}} \right| = \begin{cases} 1 & \mathbf{w}^\top \mathbf{x}^{(j)} \cdot \mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} \leq 0 \\ 0 & \text{otherwise} \end{cases}. \tag{C.58}$$

Next we are going to control $\mathbb{P}(\mathbf{w}^\top \mathbf{x}^{(j)} \cdot \mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} \leq 0)$. Denote $z_1 = \bar{\mathbf{w}}^\top \mathbf{x}^{(j)}$ and $z_2 = \bar{\mathbf{w}}_\epsilon^\top \mathbf{x}^{(j)}$, then z_1 and z_2 follow a joint Gaussian distribution and

$$\begin{aligned}
& \mathbb{P}(z_1 z_2 < 0) \\
&= \mathbb{P}(z_1 < 0, z_2 > 0) + \mathbb{P}(z_1 > 0, z_2 < 0) \\
&= \frac{2}{2\pi\sqrt{1-\rho^2}} \int_0^{+\infty} \int_{-\infty}^0 \exp\left(-\frac{1}{2(1-\rho^2)}(z_1^2 - 2\rho z_1 z_2 + z_2^2)\right) dz_1 dz_2 \\
&\leq \frac{2}{2\pi\sqrt{1-\rho^2}} \int_0^{+\infty} \int_{-\infty}^0 \exp\left(-\frac{1}{2(1-\rho^2)}(z_1^2 + z_2^2)\right) dz_1 dz_2 \\
&= \frac{2}{2\pi\sqrt{1-\rho^2}} \int_0^{+\infty} \exp\left(-\frac{z_1^2}{2(1-\rho^2)}\right) dz_1 \cdot \int_{-\infty}^0 \exp\left(-\frac{z_2^2}{2(1-\rho^2)}\right) dz_2 \\
&= \frac{2}{2\pi\sqrt{1-\rho^2}} \cdot \left(\frac{1}{2} \cdot \sqrt{2\pi(1-\rho^2)}\right)^2 \\
&= \frac{1}{2} \sqrt{1-\rho^2},
\end{aligned} \tag{C.59}$$

where $\rho = \cos \theta_\epsilon$ denotes the correlation between z_1 and z_2 , and the inequality follows since $\rho z_1 z_2 < 0$. Hence, by the definition of θ_ϵ and ρ we have

$$\mathbb{E} \left[\left(\max_{j,l} \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \left| \left(1_{\{\mathbf{w}^\top \mathbf{x}^{(j)} > 0, \mathbf{w}^\top \mathbf{x}^{(l)} > 0\}} - 1_{\{\mathbf{w}_\epsilon^\top \mathbf{x}^{(j)} > 0, \mathbf{w}_\epsilon^\top \mathbf{x}^{(l)} > 0\}} \right) \right| \right)^2 \right] \leq \sup_{\mathbf{w} \in \mathbb{S}^{m-1}} \frac{\sin \theta_\epsilon}{2}. \quad (\text{C.60})$$

And for the first term in (C.55), denote $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_K]^\top \in \mathbb{R}^{K \times m}$, we have

$$\mathbb{E} \left[\left\| \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \right\|^2 \right] = \mathbb{E} \left[\left\| \frac{1}{K^2} \mathbf{X} \mathbf{X}^\top \right\|^2 \right] = \frac{1}{K^4} \mathbb{E} [\text{s}_{\max}(\mathbf{X})^4] \quad (\text{C.61})$$

from [88, Corollary 5.35] we know that

$$\text{s}_{\max}(\mathbf{X}) \geq \sqrt{K} + \sqrt{m} + t \quad (\text{C.62})$$

holds with probability less or equal than $2 \exp\left(-\frac{t^2}{2}\right)$, in other words,

$$\mathbb{P}(\text{s}_{\max}(\mathbf{X})^4 \geq t) \leq 2 \exp\left(-\frac{\left(t^{\frac{1}{4}} - \sqrt{K} - \sqrt{m}\right)^2}{2}\right), \quad (\text{C.63})$$

and then applying the following fact, $\mathbb{E}[Z] = \int_0^\infty \mathbb{P}(Z \geq t) dt$ holds for a positive random variable, we will obtain

$$\mathbb{E}[\text{s}_{\max}(\mathbf{X})^4] = \int_0^\infty \mathbb{P}(\text{s}_{\max}(\mathbf{X})^4 \geq t) dt \leq \int_0^\infty 2 \exp\left(-\frac{\left(t^{\frac{1}{4}} - \sqrt{K} - \sqrt{m}\right)^2}{2}\right) dt, \quad (\text{C.64})$$

and by changing variable we can write

$$\begin{aligned} \mathbb{E}[\text{s}_{\max}(\mathbf{X})^4] &\leq \int_0^\infty 8 \left(t + \sqrt{K} + \sqrt{m}\right)^3 \exp\left(-\frac{t^2}{2}\right) dt \\ &= \int_0^\infty 8 \left(t^3 + \left(\sqrt{K} + \sqrt{m}\right)^3 + 3t \left(\sqrt{K} + \sqrt{m}\right)^2 + 3t^2 \left(\sqrt{K} + \sqrt{m}\right)\right) \\ &\quad \cdot \exp\left(-\frac{t^2}{2}\right) dt \\ &\leq C \cdot \left(\sqrt{K} + \sqrt{m}\right)^3. \end{aligned} \quad (\text{C.65})$$

Combining results, we can show that

$$\mathbb{E} \left[\sup_{\mathbf{w}} \|\mathbf{A}_n(\mathbf{w}) - \mathbf{A}_n(\mathbf{w}_\epsilon)\| \right] \leq \frac{(\sqrt{K} + \sqrt{m})^{\frac{3}{2}}}{K^2} \cdot \sqrt{\sup_{\mathbf{w}} \frac{\sin \theta_\epsilon}{2}}, \quad (\text{C.66})$$

where from Lemma 18 we know that $(\sin \theta_\epsilon)^2 < (\epsilon)^2$ holds when ϵ is small, hence

$$\mathbb{P} \left(\sup_{\mathbf{w}} \|\mathbf{A}_n(\mathbf{w}) - \mathbf{A}_n(\mathbf{w}_\epsilon)\| \geq \frac{t}{3} \right) \leq \frac{C}{t} \cdot \frac{(\sqrt{K} + \sqrt{m})^{\frac{3}{2}}}{K^2} \cdot \sqrt{\epsilon} \quad (\text{C.67})$$

Thus if we let

$$\exp \left(m \log \left(12 + \frac{24}{\epsilon} \right) \right) \cdot 2 \exp \left(-C \cdot n \cdot \min \{c_1 \cdot t^2, c_2 \cdot t\} \right) \leq \frac{\delta}{2}, \quad (\text{C.68})$$

$$\frac{C}{t} \cdot \frac{(\sqrt{K} + \sqrt{m})^{\frac{3}{2}}}{K^2} \cdot \sqrt{\epsilon} \leq \frac{\delta}{2}, \quad (\text{C.69})$$

and then when

$$t \geq \max \left\{ C \cdot \frac{K-1}{K} \cdot \epsilon, \frac{C \cdot \sqrt{\epsilon} \cdot (\sqrt{K} + \sqrt{m})^{\frac{3}{2}}}{\delta K^2}, \sqrt{\frac{m \log(\frac{1}{\epsilon}) - \log(\delta)}{n}} \right\}, \quad (\text{C.70})$$

set $\epsilon = \frac{\delta}{n}$, $\delta = \frac{1}{d^{-10}}$, we will obtain that with probability at least $1 - \frac{1}{d^{10}}$,

$$\|\mathbf{A}_n(\mathbf{w}) - \mathbb{E}[\mathbf{A}_n(\mathbf{w})]\| \leq \sqrt{\frac{C \cdot m \cdot \log(n)}{n}} \quad (\text{C.71})$$

holds for all $\mathbf{w} \in \mathbb{S}^{m-1}$ and remember that only the direction of \mathbf{w} matters, hence the claim also holds for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, \frac{1}{K} \|\mathbf{w}^*\|_2)$.

Together with the result that

$$\mathbb{E}[\mathbf{A}_n(\mathbf{w})] \succeq \frac{1}{2K} \cdot \mathbf{I}, \quad (\text{C.72})$$

we can conclude that when the sample size $n \geq mK^2 \cdot \log(n)$ with probability at least $1 - \frac{1}{d^{10}}$,

$$\mathbf{A}_n(\mathbf{w}) \succeq \frac{1}{4K} \cdot \mathbf{I} \quad (\text{C.73})$$

holds for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, \frac{1}{K}\|\mathbf{w}^*\|_2)$.

□

Proof of Lemma 20

Proof. Next we will show that the second term in (C.33) is uniformly upper bounded for all \mathbf{w} in a local neighborhood of \mathbf{w}^* . With some algebra we can have the following deterministic upper bound,

$$\begin{aligned} & \left| \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} \cdot \mathbf{h}^\top \mathbf{x}_i^{(l)} \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0\}} - 1_{\{\mathbf{w}^{*\top} \mathbf{x}_i^{(j)} > 0\}} \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right| \\ & \leq \left| \mathbf{h}^\top \mathbf{x}_i^{(l)} \right| \cdot \left| \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} \right| \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} \cdot \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} < 0\}}, \end{aligned} \quad (\text{C.74})$$

further since $1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} \cdot \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} < 0\}} = 1_{\{(\mathbf{w}^* + \mathbf{h})^\top \mathbf{x}_i^{(j)} \cdot \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} < 0\}} \leq 1_{\{|\mathbf{w}^{*\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}}$ holds, we will obtain the following bound,

$$\begin{aligned} \left| \mathbf{h}^\top \mathbf{x}_i^{(l)} \right| \cdot \left| \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} \right| \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} \cdot \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} < 0\}} & \leq \left| \mathbf{h}^\top \mathbf{x}_i^{(l)} \right| \cdot \left| \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} \right| \cdot 1_{\{|\mathbf{w}^{*\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}} \\ & \leq \left| \mathbf{h}^\top \mathbf{x}_i^{(l)} \right| \cdot \left| \mathbf{h}^\top \mathbf{x}_i^{(j)} \right| \cdot 1_{\{|\mathbf{w}^{*\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}}, \end{aligned}$$

i.e.,

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{w}^{*\top} \mathbf{x}_i^{(j)} \cdot \mathbf{h}^\top \mathbf{x}_i^{(l)} \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0\}} - 1_{\{\mathbf{w}^{*\top} \mathbf{x}_i^{(j)} > 0\}} \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \\ & \leq \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \left| \mathbf{h}^\top \mathbf{x}_i^{(l)} \right| \cdot \left| \mathbf{h}^\top \mathbf{x}_i^{(j)} \right| \cdot 1_{\{|\mathbf{w}^{*\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}} \\ & = \frac{1}{K^2} \sum_{j \neq l} \frac{1}{n} \sum_{i=1}^n \left| \mathbf{h}^\top \mathbf{x}_i^{(l)} \right| \cdot \left| \mathbf{h}^\top \mathbf{x}_i^{(j)} \right| \cdot 1_{\{|\mathbf{w}^{*\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}} \\ & \quad + \frac{1}{K^2} \sum_{j=1}^K \frac{1}{n} \sum_{i=1}^n \left| \mathbf{h}^\top \mathbf{x}_i^{(j)} \right|^2 \cdot 1_{\{|\mathbf{w}^{*\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}}. \end{aligned} \quad (\text{C.75})$$

The key observation here is that $\left\{ |\mathbf{w}^{*\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}| \right\}$ is an event with small probability when \mathbf{w} is close to \mathbf{w}^* , we first show the second term in (C.75) can be upper bounded by $\|\mathbf{h}\|_2^2$ up to some scaling factor and the result is summarized as follows.

Lemma 22. For any $\epsilon > 0$ and some large enough constants c_1, c_2 , if $n \geq c_1 \cdot m\epsilon^{-2} \log(\epsilon^{-1})$, then with probability at least $1 - c_2 \exp(-\epsilon^2 n)$,

$$\frac{1}{n} \sum_{i=1}^n |\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2 \cdot 1_{\{|\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}} \leq \left(\frac{1.18}{K^{\frac{3}{2}}} + \epsilon \right) \cdot \|\mathbf{h}\|_2^2, \quad (\text{C.76})$$

holds for all non-zeros $\mathbf{h} \in \mathbb{R}^m$ satisfying $\|\mathbf{h}\|_2 \leq \frac{1}{K^{\frac{3}{2}}} \|\mathbf{w}^*\|_2$, by setting $\epsilon = \frac{0.02}{K^{\frac{3}{2}}}$, we have

$$\frac{1}{n} \sum_{i=1}^n |\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2 \cdot 1_{\{|\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}} \leq \frac{1.2}{K^{\frac{3}{2}}} \cdot \|\mathbf{h}\|_2^2 \quad (\text{C.77})$$

holds with probability at least $1 - c_2 \exp(-K^3 n)$, as long as $n \geq c_1 \cdot mK^3 \log(K^{\frac{3}{2}})$ and $\|\mathbf{h}\|_2 \leq \frac{1}{K^{\frac{3}{2}}} \|\mathbf{w}^*\|_2$ hold.

Moreover, with slight adaption we can directly show a similar upper bound holds for the first term in (C.75). Hence applying Lemma 22 we have that

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} \cdot \mathbf{h}^\top \mathbf{x}_i^{(l)} \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0\}} - 1_{\{\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} > 0\}} \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \\ & \leq \frac{K^2 - K}{K^2} \cdot \frac{1.2}{K^{\frac{3}{2}}} \|\mathbf{h}\|_2^2 + \frac{K}{K^2} \cdot \frac{1.2}{K^{\frac{3}{2}}} \|\mathbf{h}\|_2^2 \leq \frac{1.2}{K^{\frac{3}{2}}} \cdot \|\mathbf{h}\|_2^2, \end{aligned} \quad (\text{C.78})$$

holds with probability at least $1 - c_2 \exp(-K^3 n)$ as long as $n \geq c_1 \cdot mK^3 \log(K^{\frac{3}{2}})$ and $\|\mathbf{h}\|_2 \leq \frac{1}{K^{\frac{3}{2}}} \|\mathbf{w}^*\|_2$ hold. \square

Proof of Lemma 21

Proof. Recall that the gradient is

$$\begin{aligned} \nabla L_n(\mathbf{w}) &= 2 \cdot \frac{1}{n} \sum_{i=1}^n \frac{1}{K^2} \sum_{j=1}^K \sum_{l=1}^K \left(\phi(\mathbf{w}^\top \mathbf{x}_i^{(j)}) - \phi(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}) \right) \phi'(\mathbf{w}^\top \mathbf{x}_i^{(l)}) \mathbf{x}_i^{(l)} \\ &= \frac{2}{K^2} \sum_{j=1}^K \sum_{l=1}^K \frac{1}{n} \sum_{i=1}^n \left(\mathbf{w}^\top \mathbf{x}_i^{(j)} \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right. \\ & \quad \left. - \mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} \cdot 1_{\{\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right) \mathbf{x}_i^{(l)} \end{aligned} \quad (\text{C.79})$$

Let $\mathbf{X}^{(l)} = [\mathbf{x}_1^{(l)}, \dots, \mathbf{x}_n^{(l)}] \in \mathbb{R}^{m \times n}$ and $\mathbf{a}^{(j,l)} = [a_1^{(j,l)}, \dots, a_n^{(j,l)}]^\top \in \mathbb{R}^n$ where $a_i^{(j,l)} = \mathbf{w}^\top \mathbf{x}_i^{(j)} \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} - \mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} \cdot 1_{\{\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}}$, then we can rewrite the gradient as

$$\nabla L_n(\mathbf{w}) = \frac{2}{K^2} \sum_{j=1}^K \sum_{l=1}^K \frac{1}{n} \mathbf{X}^{(l)} \mathbf{a}^{(j,l)} \quad (\text{C.80})$$

since the spectrum of $\mathbf{X}^{(l)}$ is well controlled when n is large enough, hence in order to upper bound the norm of $\nabla L_n(\mathbf{w})$, we will analyze the upper bound of $\|\mathbf{a}^{(j,l)}\|_2^2$. We can further rewrite $a_i^{(j,l)}$ as

$$\mathbf{h}^\top \mathbf{x}_i^{(j)} \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} + \mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0\}} - 1_{\{\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} > 0\}} \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}}$$

by triangle inequality we have that

$$\begin{aligned} & \|\mathbf{a}^{(j,l)}\|_2 \\ & \leq \sqrt{\sum_{i=1}^n \left(\mathbf{h}^\top \mathbf{x}_i^{(j)} \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0, \mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right)^2} \\ & \quad + \sqrt{\sum_{i=1}^n \left(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} \cdot \left(1_{\{\mathbf{w}^\top \mathbf{x}_i^{(j)} > 0\}} - 1_{\{\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} > 0\}} \right) \cdot 1_{\{\mathbf{w}^\top \mathbf{x}_i^{(l)} > 0\}} \right)^2} \\ & \leq \sqrt{\sum_{i=1}^n \left(\mathbf{h}^\top \mathbf{x}_i^{(j)} \right)^2} + \sqrt{\sum_{i=1}^n \left(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)} \right)^2 \cdot 1_{\{(\mathbf{w}^\top \mathbf{x}_i^{(j)}) \cdot (\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}) < 0\}}} \\ & \leq \sqrt{\sum_{i=1}^n \left(\mathbf{h}^\top \mathbf{x}_i^{(j)} \right)^2} + \sqrt{\sum_{i=1}^n \left(\mathbf{h}^\top \mathbf{x}_i^{(j)} \right)^2 \cdot 1_{\{|\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}| < |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}}} \\ & \leq 2 \sqrt{\sum_{i=1}^n \left(\mathbf{h}^\top \mathbf{x}_i^{(j)} \right)^2} \end{aligned} \quad (\text{C.81})$$

Applying Lemma 17 to the above term will give us

$$\mathbb{P} \left(\sqrt{\sum_{i=1}^n \left(\mathbf{h}^\top \mathbf{x}_i^{(j)} \right)^2} \leq \sqrt{n(1+\epsilon)} \|\mathbf{h}\|_2 \right) \geq 1 - 2 \exp(-c\epsilon^2 n), \quad (\text{C.82})$$

as long as $n \geq cm\epsilon^{-2}$. Finally we can upper bound the norm of the gradient as

$$\|\nabla L_n(\mathbf{w})\|_2^2 \leq \frac{2}{K^2} \sum_{j=1}^K \sum_{l=1}^K \left\| \frac{1}{n} \mathbf{X}^{(l)} \mathbf{a}^{(j,l)} \right\|_2^2 \leq \frac{2}{K^2} \sum_{j=1}^K \sum_{l=1}^K \frac{1}{n^2} \|\mathbf{X}^{(l)}\|^2 \cdot \|\mathbf{a}^{(j,l)}\|_2^2,$$

from [88, Corollary 5.35] we know that $\mathbb{P}(\|\mathbf{X}^{(l)}\|^2 \leq 9n) \geq 1 - 2\exp(-n/2)$ as long as $n \geq m$, hence by union bound we can conclude that

$$\mathbb{P}(\|\nabla L_n(\mathbf{w})\|_2^2 \leq 18(1+\epsilon)\|\mathbf{h}\|_2^2) \geq 1 - 2\exp(-c\epsilon^2 n) - 2\exp(-n/2) \quad (\text{C.83})$$

holds as long as $n \geq cm\epsilon^{-2}$. □

Proof of Lemma 22

Proof. We will follow the same idea in the proof of [74, Lemma 7]. Firstly, for a fixed \mathbf{h} we will apply Bernstein type concentration inequality to show the result holds, and then apply a covering argument to generalize the result. We define the following Lipschitz function,

$$F_i(t) = \begin{cases} t, & t > \left(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}\right)^2 \\ \frac{1}{\delta} \left(t - \left(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}\right)^2 \right) + \left(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}\right)^2 & (1-\delta) \cdot \left(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}\right)^2 \leq t \leq \left(\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}\right)^2 \\ 0, & \text{otherwise} \end{cases}$$

for $i = 1, \dots, n$, and it's easy to verify that

$$|\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2 \cdot \mathbb{1}_{\{|\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}| \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|\}} \leq F_i\left(|\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2\right) \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2 \cdot \mathbb{1}_{\{(1-\delta)|\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}|^2 \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2\}}.$$

Next we are going to upper bound the expectation of the right-hand side term above,

for simplicity we let $\gamma_i = \frac{|\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2}{\|\mathbf{h}\|_2^2} \cdot \mathbb{1}_{\{(1-\delta)|\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}|^2 \leq |\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2\}}$ and $r = \frac{\|\mathbf{h}\|_2}{\|\mathbf{w}^*\|_2}$, then

$$\mathbb{E}[\gamma_i] = \int \int_{-\infty}^{\infty} \mathbb{E} \left[\gamma_i |\mathbf{w}^{\star\top} \mathbf{x}_i^{(j)}| = \tau_1 \|\mathbf{w}^*\|_2, \mathbf{h}^\top \mathbf{x}_i^{(j)} = \tau_2 \|\mathbf{h}\|_2 \right] \cdot f(\tau_1, \tau_2) d\tau_1 d\tau_2, \quad (\text{C.84})$$

where $f(\tau_1, \tau_2)$ is the joint density of two Gaussian random variable with correlation $\rho = \frac{\mathbf{h}^\top \mathbf{w}^*}{\|\mathbf{h}\|_2 \|\mathbf{w}^*\|_2}$. Continuing evaluate (C.84) we will obtain

$$\begin{aligned} \mathbb{E}[\gamma_i] &= \int \int_{-\infty}^{\infty} \tau_2^2 \cdot \mathbf{1}_{\{\sqrt{1-\delta}|\tau_1| \leq |\tau_2| \cdot r\}} \cdot f(\tau_1, \tau_2) d\tau_1 d\tau_2 \\ &= \frac{1}{\sqrt{2\pi}} \int_0^{\infty} \tau_2^2 \exp\left(-\frac{\tau_2^2}{2}\right) \cdot \left(\operatorname{erf}\left(\frac{\left(\frac{r}{\sqrt{1-\delta}} - \rho\right) \tau_2}{\sqrt{1-\rho^2}}\right) + \operatorname{erf}\left(\frac{\left(\frac{r}{\sqrt{1-\delta}} + \rho\right) \tau_2}{\sqrt{1-\rho^2}}\right) \right) d\tau_2, \end{aligned}$$

we omit calculations for the second equality here, further since

$$\begin{aligned} &\operatorname{erf}\left(\frac{\left(\frac{r}{\sqrt{1-\delta}} - \rho\right) \tau_2}{\sqrt{1-\rho^2}}\right) + \operatorname{erf}\left(\frac{\left(\frac{r}{\sqrt{1-\delta}} + \rho\right) \tau_2}{\sqrt{1-\rho^2}}\right) \\ &= \operatorname{erf}\left(\frac{\left(\frac{r}{\sqrt{1-\delta}} + \rho\right) \tau_2}{\sqrt{1-\rho^2}}\right) - \operatorname{erf}\left(\frac{\left(\rho - \frac{r}{\sqrt{1-\delta}}\right) \tau_2}{\sqrt{1-\rho^2}}\right) \\ &= \left(\frac{2\frac{r}{\sqrt{1-\delta}} \cdot \tau_2}{\sqrt{1-\rho^2}}\right) \cdot \frac{2}{\sqrt{\pi}} \exp(-z^2), \end{aligned} \tag{C.85}$$

where the first equality holds since $\operatorname{erf}(\cdot)$ is an odd function, the second equality holds by applying the mean value theorem, and $z = \lambda \frac{\left(\frac{r}{\sqrt{1-\delta}} + \rho\right) \tau_2}{\sqrt{1-\rho^2}} + (1-\lambda) \frac{\left(\rho - \frac{r}{\sqrt{1-\delta}}\right) \tau_2}{\sqrt{1-\rho^2}}$ for some $\lambda \in (0, 1)$.

From Lemma 18 we know that

$$\rho > \rho^2 \geq 1 - r^2 \tag{C.86}$$

and we have that $r = \frac{\|\mathbf{h}\|_2}{\|\mathbf{w}^*\|_2} \leq \frac{1}{K^{\frac{3}{2}}}$. When the number of neurons $K \geq 2$ and δ is small, e.g. $\delta = 0.01$, then $\frac{r}{\sqrt{1-\delta}} \approx r$. Hence, $0 < \rho - \frac{r}{\sqrt{1-\delta}} < \rho + \frac{r}{\sqrt{1-\delta}}$, notice that $\exp(-z^2)$ is monotonic decreasing with respect to z , thus we have

$$\begin{aligned} \left(\frac{2\frac{r}{\sqrt{1-\delta}} \cdot \tau_2}{\sqrt{1-\rho^2}}\right) \cdot \frac{2}{\sqrt{\pi}} \exp(-z^2) &\leq 2.27 \cdot r \cdot \frac{\tau_2}{\sqrt{1-\rho^2}} \cdot \exp\left(-\left(\frac{\left(\rho - \frac{r}{\sqrt{1-\delta}}\right) \tau_2}{\sqrt{1-\rho^2}}\right)^2\right) \\ &\leq 2.27 \cdot r \cdot \frac{\tau_2}{\sqrt{1-\rho^2}} \cdot \exp\left(-\frac{0.12\tau_2^2}{1-\rho^2}\right), \end{aligned} \tag{C.87}$$

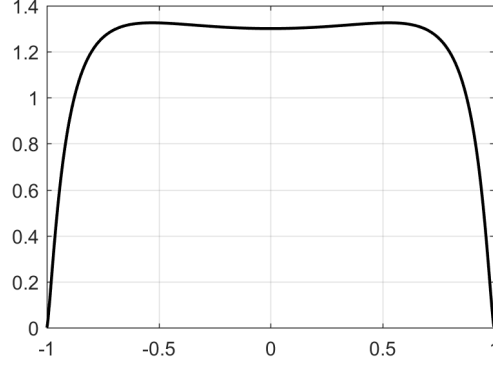


Figure C.1: Numerical integral with respect to ρ

holds, the last step follows from the fact that $\rho - \frac{r}{\sqrt{1-\delta}} \geq \sqrt{1-r^2} - \frac{r}{\sqrt{1-\delta}} \geq 0.36$, and notice that the right-hand side is monotonic decreasing, and $r \leq \frac{1}{2}$. Together we can upper bound the expectation as

$$\begin{aligned}
& \mathbb{E}[\gamma_i] \\
&= \frac{1}{\sqrt{2\pi}} \int_0^\infty \tau_2^2 \exp\left(-\frac{\tau_2^2}{2}\right) \cdot \left(\operatorname{erf}\left(\frac{\left(\frac{r}{\sqrt{1-\delta}} - \rho\right) \tau_2}{\sqrt{1-\rho^2}}\right) + \operatorname{erf}\left(\frac{\left(\frac{r}{\sqrt{1-\delta}} + \rho\right) \tau_2}{\sqrt{1-\rho^2}}\right) \right) d\tau_2 \\
&\leq \frac{2.27}{\sqrt{2\pi}} \cdot r \cdot \int_0^\infty \tau_2^3 \exp\left(-\frac{\tau_2^2}{2}\right) \cdot \frac{1}{\sqrt{1-\rho^2}} \exp\left(-\frac{0.12\tau_2^2}{1-\rho^2}\right) d\tau_2, \\
&\leq \frac{1.18}{K^{\frac{3}{2}}},
\end{aligned}$$

where we evaluate the integral in the first inequality over different choices of ρ in Fig C.1, we see that the integral is upper bounded by a constant, hence we can obtain the last inequality. Moreover, $\|F_i(|\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2)\|_{\psi_1} \leq C \cdot \|\mathbf{h}\|_2^2$ for all $i = 1, \dots, n$, hence when $r = \frac{\|\mathbf{h}\|_2}{\|\mathbf{w}^*\|_2} \leq \frac{1}{K}$ we can have

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \frac{F_i(|\mathbf{h}^\top \mathbf{x}_i^{(j)}|^2)}{\|\mathbf{h}\|_2^2} \geq \left(\frac{1.18}{K^{\frac{3}{2}}} + \epsilon\right)\right) < \exp(-c \cdot n \cdot \epsilon^2), \quad (\text{C.88})$$

holds for some constant c . Since the covering argument will be the same as the proof of [74, Lemma 7], we omit the repetition and conclude.

□

Bibliography

- [1] P. T. Boufounos and R. G. Baraniuk, “1-bit compressive sensing,” in *Information Sciences and Systems, 2008. CISS 2008. 42nd Annual Conference on*. IEEE, 2008, pp. 16–21.
- [2] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, “The convex geometry of linear inverse problems,” *Foundations of Computational Mathematics*, vol. 12, no. 6, pp. 805–849, 2012.
- [3] E. J. Candès and B. Recht, “Exact matrix completion via convex optimization,” *Foundations of Computational mathematics*, vol. 9, no. 6, p. 717, 2009.
- [4] E. J. Candès and T. Tao, “The power of convex relaxation: Near-optimal matrix completion,” *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2053–2080, 2010.
- [5] E. J. Candès, X. Li, and M. Soltanolkotabi, “Phase retrieval via wirtinger flow: Theory and algorithms,” *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1985–2007, 2015.
- [6] Y. Chen and E. Candes, “Solving random quadratic systems of equations is nearly as easy as solving linear systems,” in *Advances in Neural Information Processing Systems*, 2015, pp. 739–747.
- [7] K. Zhong, Z. Song, P. Jain, P. L. Bartlett, and I. S. Dhillon, “Recovery guarantees for one-hidden-layer neural networks,” in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70, 2017, pp. 4140–4149.
- [8] S. Lloyd, “Least squares quantization in PCM,” *IEEE transactions on information theory*, vol. 28, no. 2, pp. 129–137, 1982.
- [9] D. Donoho, “Compressed sensing,” *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, April 2006.
- [10] E. J. Candès and T. Tao, “Near-optimal signal recovery from random projections: Universal encoding strategies?” *IEEE Transactions on Information Theory*, vol. 52, no. 12, pp. 5406–5425, Dec. 2006.

- [11] J. A. Tropp, J. N. Laska, M. F. Duarte, J. K. Romberg, and R. G. Baraniuk, “Beyond Nyquist: Efficient sampling of sparse bandlimited signals,” *Information Theory, IEEE Transactions on*, vol. 56, no. 1, pp. 520–544, 2010.
- [12] M. Mishali, Y. C. Eldar, and A. J. Elron, “Xampling: Signal acquisition and processing in union of subspaces,” *IEEE Transactions on Signal Processing*, vol. 59, no. 10, pp. 4719–4734, 2011.
- [13] M. Wakin, S. Becker, E. Nakamura, M. Grant, E. Sovero, D. Ching, J. Yoo, J. Romberg, A. Emami-Neyestanak, and E. Candes, “A nonuniform sampler for wideband spectrally-sparse environments,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 2, no. 3, pp. 516–529, 2012.
- [14] F. Zeng, C. Li, and Z. Tian, “Distributed compressive spectrum sensing in cooperative multihop cognitive networks,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 1, pp. 37–48, 2011.
- [15] Y. L. Polo, Y. Wang, A. Pandharipande, and G. Leus, “Compressive wide-band spectrum sensing,” in *Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on*. IEEE, 2009, pp. 2337–2340.
- [16] O. Mehanna and N. Sidiropoulos, “Frugal sensing: Wideband power spectrum sensing from few bits,” *Signal Processing, IEEE Transactions on*, vol. 61, no. 10, pp. 2693–2703, May 2013.
- [17] Y. Chi and H. Fu, “Subspace learning from bits,” *IEEE Transactions on Signal Processing*, vol. 65, no. 17, pp. 4429–4442, Sept 2017.
- [18] R. H. Walden, “Analog-to-digital converter survey and analysis,” *IEEE Journal on selected areas in communications*, vol. 17, no. 4, pp. 539–550, 1999.
- [19] J. N. Laska, P. T. Boufounos, M. A. Davenport, and R. G. Baraniuk, “Democracy in action: Quantization, saturation, and compressive sensing,” *Applied and Computational Harmonic Analysis*, vol. 31, no. 3, pp. 429–443, 2011.
- [20] M. Slawski and P. Li, “Linear signal recovery from b -bit-quantized linear measurements: precise analysis of the trade-off between bit depth and number of measurements,” *arXiv preprint arXiv:1607.02649*, 2016.
- [21] H. L. Van Trees, *Detection, estimation, and modulation theory*. John Wiley & Sons, 2004.
- [22] Y. Chi, L. Scharf, A. Pezeshki, and A. Calderbank, “Sensitivity to basis mismatch in compressed sensing,” *IEEE Transactions on Signal Processing*, vol. 59, no. 5, pp. 2182–2195, May 2011.

- [23] G. Tang, B. N. Bhaskar, P. Shah, and B. Recht, “Compressed sensing off the grid,” *IEEE transactions on information theory*, vol. 59, no. 11, pp. 7465–7490, 2013.
- [24] Y. Chi and Y. Chen, “Compressive two-dimensional harmonic retrieval via atomic norm minimization,” *IEEE Transactions on Signal Processing*, vol. 63, no. 4, pp. 1030–1042, 2015.
- [25] Y. Li and Y. Chi, “Off-the-grid line spectrum denoising and estimation with multiple measurement vectors,” *Signal Processing, IEEE Transactions on*, vol. 64, no. 5, pp. 1257–1269.
- [26] R. Heckel and M. Soltanolkotabi, “Generalized line spectral estimation via convex optimization,” *IEEE Transactions on Information Theory*, vol. 64, no. 6, pp. 4001–4023, 2017.
- [27] Z. Yang and L. Xie, “Exact joint sparse frequency recovery via optimization methods,” *IEEE Transactions on Signal Processing*, vol. 64, no. 19, pp. 5145–5157, 2016.
- [28] G. Tang, B. N. Bhaskar, and B. Recht, “Near minimax line spectral estimation,” *IEEE Transactions on Information Theory*, vol. 61, no. 1, pp. 499–512, 2015.
- [29] B. N. Bhaskar, G. Tang, and B. Recht, “Atomic norm denoising with applications to line spectral estimation,” *IEEE Transactions on Signal Processing*, vol. 61, no. 23, pp. 5987–5999, 2013.
- [30] Y. Chi and M. F. Da Costa, “Harnessing sparsity over the continuum: Atomic norm minimization for super resolution,” *arXiv preprint arXiv:1904.04283*, 2019.
- [31] G. Gray and G. Zeoli, “Quantization and saturation noise due to analog-to-digital conversion,” *IEEE Transactions on Aerospace and Electronic Systems*, no. 1, pp. 222–223, 1971.
- [32] A. Gupta, R. Nowak, and B. Recht, “Sample complexity for 1-bit compressed sensing and sparse classification,” in *Information Theory Proceedings (ISIT), 2010 IEEE International Symposium on*. IEEE, 2010, pp. 1553–1557.
- [33] M. Yan, Y. Yang, and S. Osher, “Robust 1-bit compressive sensing using adaptive outlier pursuit,” *Signal Processing, IEEE Transactions on*, vol. 60, no. 7, pp. 3868–3875, 2012.
- [34] Y. Plan and R. Vershynin, “One-bit compressed sensing by linear programming,” *Communications on Pure and Applied Mathematics*, vol. 66, no. 8, pp. 1275–1297, 2013.

- [35] —, “Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach,” *IEEE Transactions on Information Theory*, vol. 59, no. 1, pp. 482–494, 2013.
- [36] —, “The generalized Lasso with non-linear observations,” *IEEE Transactions on Information Theory*, vol. 62, no. 3, pp. 1528–1537, 2016.
- [37] Y. Plan, R. Vershynin, and E. Yudovina, “High-dimensional estimation with geometric constraints,” *Information and Inference: A Journal of the IMA*, vol. 6, no. 1, pp. 1–40, 2016.
- [38] L. Jacques, J. N. Laska, P. T. Boufounos, and R. G. Baraniuk, “Robust 1-bit compressive sensing via binary stable embeddings of sparse vectors,” *Information Theory, IEEE Transactions on*, vol. 59, no. 4, pp. 2082–2102, 2013.
- [39] S. Oymak and M. Soltanolkotabi, “Fast and reliable parameter estimation from nonlinear observations,” *SIAM Journal on Optimization*, vol. 27, no. 4, pp. 2276–2300, 2017.
- [40] P. T. Boufounos, L. Jacques, F. Krahmer, and R. Saab, “Quantization and compressive sensing,” in *Compressed Sensing and its Applications*. Springer, 2015, pp. 193–237.
- [41] J. N. Laska and R. G. Baraniuk, “Regime change: Bit-depth versus measurement-rate in compressive sensing,” *IEEE Transactions on Signal Processing*, vol. 60, no. 7, pp. 3496–3505, 2012.
- [42] A. Host-Madsen and P. Handel, “Effects of sampling and quantization on single-tone frequency estimation,” *IEEE Transactions on Signal Processing*, vol. 48, no. 3, pp. 650–662, 2000.
- [43] P. Pakrooh, L. L. Scharf, A. Pezeshki, and Y. Chi, “Analysis of fisher information and the Cramér-Rao bound for nonlinear parameter estimation after compressed sensing,” in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 6630–6634.
- [44] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [45] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

- [46] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [47] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, and M. Lanctot, “Mastering the game of Go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [48] M. Soltanolkotabi, “Learning ReLUs via gradient descent,” in *Advances in Neural Information Processing Systems 30*, 2017, pp. 2007–2017.
- [49] Y. Li and Y. Yuan, “Convergence analysis of two-layer neural networks with ReLU activation,” in *Advances in Neural Information Processing Systems*, 2017, pp. 597–607.
- [50] S. Mei, Y. Bai, and A. Montanari, “The landscape of empirical risk for nonconvex losses,” *The Annals of Statistics*, vol. 46, no. 6A, pp. 2747–2774, 2018.
- [51] R. Sun and Z.-Q. Luo, “Guaranteed matrix completion via non-convex factorization,” *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 6535–6579, 2016.
- [52] Y. Chen and Y. Chi, “Harnessing structures in big data via guaranteed low-rank matrix estimation,” *IEEE Signal Processing Magazine*, vol. 35, no. 4, pp. 14–31, 2018.
- [53] E. J. Candès, X. Li, and M. Soltanolkotabi, “Phase retrieval via Wirtinger flow: Theory and algorithms,” *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1985–2007, April 2015.
- [54] R. Ge and T. Ma, “On the optimization landscape of tensor decompositions,” in *Advances in Neural Information Processing Systems 30*, 2017, pp. 3653–3663.
- [55] R. Ge, J. D. Lee, and T. Ma, “Matrix completion has no spurious local minimum,” in *Advances in Neural Information Processing Systems*, 2016, pp. 2973–2981.
- [56] J. Sun, Q. Qu, and J. Wright, “Complete dictionary recovery using nonconvex optimization,” *International Conference on Machine Learning*, pp. 2351–2360, 2015.
- [57] S. Bhojanapalli, B. Neyshabur, and N. Srebro, “Global optimality of local search for low rank matrix recovery,” in *Advances in Neural Information Processing Systems*, 2016, pp. 3873–3881.

- [58] C. Ma, K. Wang, Y. Chi, and Y. Chen, “Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion,” in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 3345–3354.
- [59] Y. Chi, Y. M. Lu, and Y. Chen, “Nonconvex optimization meets low-rank matrix factorization: An overview,” *IEEE Transactions on Signal Processing*, vol. 67, no. 20, pp. 5239–5269, 2019.
- [60] M. Soltanolkotabi, A. Javanmard, and J. D. Lee, “Theoretical insights into the optimization landscape of over-parameterized shallow neural networks,” *IEEE Transactions on Information Theory*, 2018.
- [61] D. Boob and G. Lan, “Theoretical properties of the global optimizer of two layer neural network,” *arXiv preprint arXiv:1710.11241*, 2017.
- [62] I. Safran and O. Shamir, “On the quality of the initial basin in overspecified neural networks,” in *International Conference on Machine Learning*, 2016, pp. 774–782.
- [63] Q. Nguyen and M. Hein, “The loss surface of deep and wide neural networks,” in *International Conference on Machine Learning*, 2017, pp. 2603–2612.
- [64] Y. Tian, “An analytical formula of population gradient for two-layered ReLU network and its applications in convergence and critical point analysis,” in *Proceedings of the 34th International Conference on Machine Learning*, 2017, pp. 3404–3413.
- [65] R. Ge, J. D. Lee, and T. Ma, “Learning one-hidden-layer neural networks with landscape design,” in *International Conference on Learning Representations*, 2018.
- [66] I. Safran and O. Shamir, “Spurious local minima are common in two-layer relu neural networks,” *arXiv preprint arXiv:1712.08968*, 2017.
- [67] Y. Chen, Y. Chi, J. Fan, and C. Ma, “Gradient descent with random initialization: fast global convergence for nonconvex phase retrieval,” *Mathematical Programming*, vol. 176, no. 1-2, pp. 5–37, 2019.
- [68] S. Oymak, “Learning compact neural networks with regularization,” in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 3966–3975.
- [69] X. Zhang, Y. Yu, L. Wang, and Q. Gu, “Learning One-hidden-layer ReLU networks via Gradient Descent,” in *Proceedings of Machine Learning Research*, ser.

- Proceedings of Machine Learning Research, K. Chaudhuri and M. Sugiyama, Eds., vol. 89. PMLR, 16–18 Apr 2019, pp. 1524–1534.
- [70] Y. Li, C. Ma, Y. Chen, and Y. Chi, “Nonconvex matrix factorization from rank-one measurements,” in *Proceedings of Machine Learning Research*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and M. Sugiyama, Eds., vol. 89. PMLR, 16–18 Apr 2019, pp. 1496–1505.
- [71] A. Brutzkus and A. Globerson, “Globally optimal gradient descent for a ConvNet with Gaussian inputs,” in *Proceedings of the 34th International Conference on Machine Learning*, 2017, pp. 605–614.
- [72] S. S. Du, J. D. Lee, and Y. Tian, “When is a convolutional filter easy to learn?” in *International Conference on Learning Representations*, 2018.
- [73] S. Du, J. Lee, Y. Tian, A. Singh, and B. Póczos, “Gradient descent learns one-hidden-layer CNN: Dont be afraid of spurious local minima,” in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 1339–1348.
- [74] H. Zhang, Y. Zhou, Y. Liang, and Y. Chi, “A nonconvex approach for phase retrieval: Reshaped wirtinger flow and incremental algorithms,” *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 5164–5198, 2017.
- [75] T. Adali, P. J. Schreier, and L. L. Scharf, “Complex-valued signal processing: The proper way to deal with impropriety,” *IEEE Transactions on Signal Processing*, vol. 59, no. 11, pp. 5101–5125, 2011.
- [76] M. Grant and S. Boyd, “CVX: Matlab software for disciplined convex programming.”
- [77] N. Boyd, G. Schiebinger, and B. Recht, “The alternating descent conditional gradient method for sparse inverse problems,” *SIAM Journal on Optimization*, vol. 27, no. 2, pp. 616–639, 2017.
- [78] N. Rao, P. Shah, and S. Wright, “Forward–backward greedy algorithms for atomic norm regularization,” *IEEE Transactions on Signal Processing*, vol. 63, no. 21, pp. 5798–5811, 2015.
- [79] K. Zhong, Z. Song, and I. S. Dhillon, “Learning non-overlapping convolutional neural networks with multiple kernels,” *arXiv preprint arXiv:1711.03440*, 2017.
- [80] H. Fu and Y. Chi, “Quantized spectral compressed sensing: Cramér-Rao bounds and recovery algorithms,” *IEEE Transactions on Signal Processing*, vol. 66, no. 12, pp. 3268–3279, 2018.

- [81] —, “Principal subspace estimation for low-rank toeplitz covariance matrices with binary sensing,” in *2016 50th Asilomar Conference on Signals, Systems and Computers*. IEEE, 2016, pp. 1344–1348.
- [82] H. Fu, Y. Chi, and Y. Liang, “Guaranteed recovery of one-hidden-layer neural networks via cross entropy,” *arXiv preprint arXiv:1802.06463*, 2018.
- [83] Y. Chen, Y. Chi, and A. J. Goldsmith, “Robust and universal covariance estimation from quadratic measurements via convex programming,” in *Information Theory (ISIT), 2014 IEEE International Symposium on*. IEEE, 2014, pp. 2017–2021.
- [84] J.-F. Cai, X. Qu, W. Xu, and G.-B. Ye, “Robust recovery of complex exponential signals from random gaussian projections via low rank hankel matrix reconstruction,” *Applied and computational harmonic analysis*, vol. 41, no. 2, pp. 470–490, 2016.
- [85] R. Vershynin, “Introduction to the non-asymptotic analysis of random matrices,” *Compressed Sensing, Theory and Applications*, pp. 210 – 268, 2012.
- [86] E. J. Candès, T. Strohmer, and V. Voroninski, “Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming,” *Communications on Pure and Applied Mathematics*, vol. 66, no. 8, pp. 1241–1274, 2013.
- [87] W. Rudin *et al.*, *Principles of mathematical analysis*. McGraw-hill New York, 1964, vol. 3.
- [88] R. Vershynin, *Introduction to the non-asymptotic analysis of random matrices*. Cambridge University Press, 2012, p. 210268.