

# Tangent Estimation from Point Samples\*

Siu-Wing Cheng<sup>†</sup> Man-Kwun Chiu<sup>‡§</sup>

## Abstract

Let  $\mathcal{M}$  be an  $m$ -dimensional smooth compact manifold embedded in  $\mathbb{R}^d$ , where  $m$  is a constant known to us. Suppose that a dense set of points are sampled from  $\mathcal{M}$  according to a Poisson process with an unknown parameter. Let  $p$  be any sample point, let  $\varrho$  be the local feature size at  $p$ , and let  $\varrho\varepsilon$  be the distance from  $p$  to the  $(n+1)$ th nearest sample point for some  $n$  between  $\binom{m+1}{2} + 1$  and  $\binom{d+1}{2}$ . Using the  $n$  sample points nearest to  $p$ , we can estimate the tangent space at  $p$  and it holds with probability  $1 - O(n^{-1/3})$  that the angular error is  $O(\varepsilon^2)$ . The running time is bounded by the time to compute the thin SVD of an  $n \times \binom{d+1}{2}$  matrix and the full SVD of an  $n \times d$  matrix, which is usually  $O(d^2n^2)$  in practice. We implemented the algorithm and experimentally verified its effectiveness on both noiseless and noisy data.

---

\*Research supported by the Research Grant Council, Hong Kong, China (project no. 612109). Part of the work was done while Chiu was at HKUST.

<sup>†</sup>Dept. of Computer Science and Engineering, HKUST, Hong Kong.

<sup>‡</sup>National Institute of Informatics (NII), Tokyo, Japan.

<sup>§</sup>JST, ERATO, Kawarabayashi Large Graph Project.

# 1 Introduction

Data points corresponding to experimental observations commonly reside in  $\mathbb{R}^d$  for some large  $d$ , but it is often postulated that the data points live on an unknown manifold  $\mathcal{M}$  of much lower dimension. Indeed, manifold learning has been applied in various problems such as network anomaly detection, image segmentation, and object tracking in video [6, 29]. The goal is to learn the manifold structure from sample points, including the intrinsic dimension, tangent spaces, and ultimately a faithful reconstruction. Theoretical algorithms have been developed to obtain faithful reconstructions [3, 4, 9], but their practical performance is unclear. We reexamine the key tasks in the problem to put our result in context.

The first task is to detect the manifold dimension. Many effective methods have already been developed in machine learning (e.g. [21, 22, 25, 27, 33]) and computational geometry [8, 10, 11, 15].

The second task is to estimate the tangent space at the sample points. Approximate tangent spaces at the sample points are needed in [9] to compute the cocone complex, which becomes a faithful reconstruction after removing slivers. Approximate tangent spaces at the sample points are also needed in [3] in order to form the tangential Delaunay complex, from which a faithful reconstruction is extracted after sliver removal. Tangent space estimation also finds application in clustering data points from multiple manifolds that may intersect each other [14, 18]. The tangent space estimation at a sample point  $\mathbf{p}$  has been explicitly or implicitly performed in many previous work by fitting an affine subspace to the sample points in a neighborhood of  $\mathbf{p}$  [1, 2, 10, 15, 26, 30, 31, 33]. An alternative method is based on analyzing the Voronoi cell of  $\mathbf{p}$  [11]. The error measure for tangent space estimation is the angular error, which is the maximum angle between a vector in the true tangent space at  $\mathbf{p}$  and the projection of that vector in the estimated tangent space. Bounds on the angular error have been proved (explicitly or implicitly) for the methods in [2, 10, 11, 15, 26], and the emphasis is on how these bounds depend on the sampling density. The radii of the neighborhoods used by the methods in [2, 10, 11, 15] for tangent estimation decrease as the sampling density increases, and their angular error bounds are *linear* in the ratio of the neighborhood radius to the local feature size at  $\mathbf{p}$ .

Let  $m$  be the dimension of the unknown manifold  $\mathcal{M}$ . We present a method to estimate the tangent space at a sample point  $\mathbf{p}$  using the sample points in a local neighborhood. Let  $\varrho$  be the local feature size of  $\mathcal{M}$  at  $\mathbf{p}$ . Let  $r$  be the neighborhood radius so that all sample points at distance less than  $r$  from  $\mathbf{p}$  are used in the tangent space estimation. Our method gives an angular error bound of  $O((r/\varrho)^2)$  radians with high probability, provided that  $\mathcal{M}$  is smooth and compact, the sample points are drawn from  $\mathcal{M}$  according to a Poisson process with an unknown parameter, and the manifold dimension  $m$  is a constant known to us.<sup>1</sup> Within a neighborhood of radius  $r$  from  $\mathbf{p}$ , the tangent space can rotate by at most  $O(r/\varrho)$  radians [9], where  $\varrho$  is the local feature size at  $\mathbf{p}$ , and the turning angle may sometimes be  $\Omega(r/\varrho)$ ; for example, when  $\mathcal{M}$  is the unit sphere  $\mathbb{S}^m$ . An angular error bound linear in  $r/\varrho$  is thus asymptotically as good as taking the tangent space at an arbitrary sample point in the neighborhood as the approximation. We do better as our angular error bound is  $O((r/\varrho)^2)$ . We elaborate on our result and compare it with previous works that provide angular error bounds in Section 1.2.

---

<sup>1</sup>Since our method reduces to solving an eigenvalue problem, an appropriate thresholding of the eigenvalues should determine  $m$ . We do not pursue automatic dimension detection in this article in order to focus on the tangent estimation. We comment on the determination of  $m$  further in the conclusion.

## 1.1 Notation

An uppercase letter in `mathsf` font denotes a matrix and the corresponding italic lowercase letter with subscripts denotes the matrix elements. For example,  $Z$  denotes a matrix;  $z_{ij}$  denotes the  $(i, j)$  entry of  $Z$ ;  $\mathbf{z}_{i*}$  and  $\mathbf{z}_{*j}$  denote the  $i$ th row vector and the  $j$ th column vector of  $Z$ , respectively. Similarly,  $\mathbf{v}$  denotes a vector and  $v_i$  denotes the  $i$ th coordinate of  $\mathbf{v}$ .  $Z^t$  and  $\mathbf{v}^t$  denote the transposes. We use  $I_j$  to denote a  $j \times j$  identity matrix,  $\mathbf{0}_{i,j}$  an  $i \times j$  zero matrix, and  $\text{diag}_j(\sigma_1, \sigma_2, \dots, \sigma_j)$  a  $j \times j$  diagonal matrix with entries  $\sigma_1, \dots, \sigma_j$  in this order. We reserve  $\mathbf{0}$  to denote the origin of  $\mathbb{R}^d$ .

The 2-norms of  $\mathbf{v}$  and  $Z$  are  $\|\mathbf{v}\| = (\sum_i v_i^2)^{1/2}$  and  $\|Z\| = \max\{\|Z\mathbf{v}\| : \|\mathbf{v}\| = 1\}$ . If  $Z$  is symmetric,  $\|Z\|$  also equals  $\max\{|\mathbf{v}^t Z \mathbf{v}| : \|\mathbf{v}\| = 1\}$ . The Frobenius norm of  $Z$  is  $\|Z\|_F = (\sum_i \sum_j z_{ij}^2)^{1/2}$ . It is known that  $\|Z\| \leq \|Z\|_F \leq \sqrt{k} \|Z\|$ , where  $k$  is the number of rows or columns in  $Z$ , whichever is smaller [17].

Given a square matrix  $Z$ , a vector  $\mathbf{v}$  is an eigenvector of  $Z$  if and only if  $Z\mathbf{v} = \lambda\mathbf{v}$  for some  $\lambda \in \mathbb{R}$ , and  $\lambda$  is known as an eigenvalue of  $Z$ . If  $Z$  has dimension  $k$ , then  $Z$  has at most  $k$  real eigenvalues. If  $Z$  is symmetric as well, it has  $k$  real eigenvalues.

The *thin* singular value decomposition (thin SVD) of a  $k \times l$  matrix  $Z$ ,  $k \leq l$ , is a product  $LDR^t$ , where  $L$  is a  $k \times k$  matrix consisting of unit eigenvectors of  $ZZ^t$ ,  $D$  is a  $k \times k$  diagonal matrix consisting of the singular values of  $Z$  (i.e., square roots of the eigenvalues of  $ZZ^t$ ), and  $R$  is an  $l \times k$  matrix formed by  $k$  of the unit eigenvectors of  $Z^t Z$  corresponding to the  $k$  largest eigenvalues. We assume that the singular values of  $Z$  are in descending order on the diagonal of  $D$ . The *full* SVD of the same matrix  $Z$  is  $L(D \ \mathbf{0}_{k,l-k})\bar{R}^t$ , where  $\bar{R}$  is an  $l \times l$  matrix formed by the  $l$  unit eigenvectors of  $Z^t Z$  and  $R$  is the leftmost  $l \times k$  submatrix of  $\bar{R}$ .

Given a diagonal square matrix  $D$ , its *pseudoinverse*  $D^\dagger$  is obtained by replacing each non-zero entry by its reciprocal and leaving the zero entries in place. The pseudoinverse of a general matrix  $Z$  with thin SVD  $LDR^t$  and full SVD  $L(D \ \mathbf{0}_{k,l-k})\bar{R}^t$  is  $Z^\dagger = RD^\dagger L^t = \bar{R}(D^\dagger \ \mathbf{0}_{k,l-k})^t L^t$ . When  $Z$  is square and invertible,  $Z^\dagger$  is just  $Z^{-1}$ .

The largest singular value of  $Z$  is equal to  $\|Z\|$ . The positive singular values of  $Z^\dagger$  are the reciprocals of the positive singular values of  $Z$ . Therefore,  $\|Z^\dagger\|$  is the reciprocal of the smallest positive singular value of  $Z$ .

Let  $x_1, x_2, \dots, x_d$  be a fixed set of orthogonal axes throughout this paper, forming the default coordinate system of  $\mathbb{R}^d$ . The coordinates of the input sample points are expressed with respect to this coordinate system.

We are given a set of sample points drawn from  $\mathcal{M}$  according to a Poisson process with an unknown parameter  $\lambda$ : (i) for any compact subset  $B$  of  $\mathcal{M}$ , the probability that there are  $k$  points in  $B$  is  $\frac{\lambda^k \text{vol}(B)^k}{k!} e^{-\lambda \text{vol}(B)}$ , and (ii) for any disjoint compact subsets  $B_1, \dots, B_j$  of  $\mathcal{M}$ , the probability that there are  $k_i$  points in  $B_i$  for  $i \in [1, j]$  is  $\prod_{i=1}^j \frac{\lambda^{k_i} \text{vol}(B_i)^{k_i}}{(k_i)!} e^{-\lambda \text{vol}(B_i)}$ . Given such a Poisson process and on the condition that there are  $k$  sample points in a compact subset  $B \subset \mathcal{M}$ , these  $k$  sample points are uniformly distributed in  $B$  [5].

By translation, we assume without loss of generality that the origin is a sample point. Let  $\mathcal{T}$  denote the tangent space of  $\mathcal{M}$  at the origin, which is an  $m$ -dimensional vector space in  $\mathbb{R}^d$ . Every vector in  $\mathcal{T}$  has  $d$  coordinates although  $\mathcal{T}$  has dimension  $m$ . The medial axis of  $\mathcal{M}$  is the closure of the set of points in  $\mathbb{R}^d$  that have two or more closest points in  $\mathcal{M}$ . The local feature size of a point in  $\mathcal{M}$  is the distance from that point to the medial axis. Let  $\varrho$  denote the local feature size of  $\mathcal{M}$  at the origin. Let  $\{\mathbf{a}_p : p \in [1, n]\}$  denote the  $n$  sample points nearest to the origin. Let  $\varrho\varepsilon$  denote the distance from the origin to the  $(n+1)$ -th nearest sample point, where  $\varepsilon \in (0, 1)$  and  $\varepsilon$  decreases as the sampling density increases.

The manifold dimension  $m$  is treated as a constant. So we often absorb a function of  $m$  into the hidden constants in the big-Oh, big-Theta and big-Omega notation. We keep these

hidden constants scale independent; for example, the dependence on  $\varrho$  is explicitly stated. The ambient space dimension  $d$  is not a constant because one can embed  $\mathcal{M}$  in an Euclidean space of arbitrarily high dimension.

## 1.2 Main result and comparison with previous work

The intuition behind our strategy is to compute a smooth approximation of  $\mathcal{M}$  locally around the origin. Let  $\gamma_1, \dots, \gamma_d$  be any  $d$  orthogonal coordinate axes of  $\mathbb{R}^d$  such that  $\gamma_1, \dots, \gamma_m$  span  $\mathcal{T}$ . Let  $\psi$  be an (unknown) orthonormal transformation such that for every point  $\mathbf{y} \in \mathcal{M}$ ,  $\psi(\mathbf{y})$  are the coordinates of  $\mathbf{y}$  with respect to the coordinate system  $(\gamma_1, \dots, \gamma_d)$ . By the implicit function theorem, for every point  $\mathbf{y} \in \mathcal{M}$  close enough to the origin and every  $\ell \in [m+1, d]$ , the  $\ell$ -th coordinate of  $\psi(\mathbf{y})$  can be expressed as a function  $f_\ell : \mathbb{R}^m \rightarrow \mathbb{R}$  in the first  $m$  coordinates of  $\psi(\mathbf{y})$ . We call  $\{f_\ell : \ell \in [m+1, d]\}$  the *coordinate functions* of  $\mathcal{M}$  at the origin with respect to  $(\gamma_1, \dots, \gamma_d)$ .

We will approximate  $f_\ell$  by an ‘‘almost quadratic’’ function  $\widehat{F}_\ell : \mathbb{R}^d \rightarrow \mathbb{R}$  via solving an eigenvalue problem. There is not enough data to define the  $\widehat{F}_\ell$ ’s unambiguously because there are only  $n \leq \binom{d+1}{2}$  sample points. A popular approach is to add a penalty function, but a penalty function usually involves some parameter(s) and it is unclear how to tune them to obtain guarantees on the angular error. This parameter tuning phase may also be time-consuming. We also use a penalty function that involves a positive parameter. Our innovation is pushing this parameter to zero in the limit and obtain a modified eigenvalue problem. Hence, no parameter needs to be tuned and no training is required in the end. Solving this modified eigenvalue problem is equivalent to minimizing a measure of ‘‘curviness’’ of the fitting solution, which implies a theoretical guarantee on the angular error.

Our main result is stated in the following theorem. Let  $T_{\text{tsvd}}(i, j)$  and  $T_{\text{fsvd}}(i, j)$  denote the time to construct the thin and full singular value decompositions of an  $i \times j$  matrix, respectively.

**Theorem 1.1** *Suppose that  $\mathcal{M}$  is a smooth compact  $m$ -dimensional manifold in  $\mathbb{R}^d$ , where  $m$  is a constant known to us, and that points are sampled from  $\mathcal{M}$  according to a Poisson process with an unknown parameter. Assume that the origin is a sample point and its nearest  $n$  sample points are given, where  $\binom{m+1}{2} + 1 \leq n \leq \binom{d+1}{2}$ . Let  $\varrho\varepsilon$  be the distance from the origin to the  $(n+1)$ -th nearest sample point, where  $\varrho$  is the local feature size of  $\mathcal{M}$  at the origin and  $\varepsilon$  is a value in  $(0, 1)$ . We can compute in  $O(T_{\text{tsvd}}(n, \binom{d+1}{2}) + T_{\text{fsvd}}(n, d) + d^2n)$  time  $m$  orthogonal coordinate axes that span the approximate tangent space at the origin. If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , the angular error is  $O(\varepsilon^2)$ .*

The running time  $O(T_{\text{tsvd}}(n, \binom{d+1}{2}) + T_{\text{fsvd}}(n, d) + d^2n)$  of our tangent estimation algorithm is  $O(d^2n^2)$  in practice [7, 16, 17]. The worst-case running time is asymptotically bounded by the worst-case running time of multiplying an  $n \times n^r$  matrix with an  $n^r \times n$  matrix for some  $r$ . The exact bound has a sophisticated expression depending on  $r$  [24]. As two examples in our case, if  $d = O(n)$ , the time bound is  $O(n^{3.256689})$ , and if  $d = O(n^2)$ , the time bound is  $O(n^{5.180715})$ . In general, the worst-case running time is slightly better than  $O(d^2n^2)$ . Although the local feature size  $\varrho$  at the origin is used to obtain  $\varepsilon$  for expressing the angular error bound, our algorithm does not need to know  $\varrho$ .

In addition to developing an algorithm for the tangent estimation problem, we also develop some useful results along the way that may be of independent interests. Taubin gave a method for converting the curve reconstruction problem to an eigenvalue problem when there are enough sample points. We generalize this method for manifold reconstruction in high dimensions when there are insufficient sample points (Section 4). We also derive some concentration bounds on sums of powers of the coordinates of the sample points, which may be useful for other statistical analysis (Section 5).

Consider the condition  $\binom{m+1}{2} + 1 \leq n \leq \binom{d+1}{2}$  in Theorem 1.1. The formulation of our approach in Section 2 requires that  $n \geq \binom{m+1}{2} + 1$ . Notice that  $\binom{m+1}{2} + 1 < \binom{d+1}{2}$  for  $m \leq d-1$ . Our techniques are not designed for the case of  $n > \binom{d+1}{2}$ . In a manifold learning context, it is predominantly the case that  $d$  is large. Therefore, the requirement of  $n \leq \binom{d+1}{2}$  is not an issue because it is very likely that there are fewer than  $\binom{d+1}{2}$  sample points nearby, and even if there are so many sample points nearby, it is computationally less efficient to use them all. If  $d$  is not large and  $n > \binom{d+1}{2}$ , our result can still be applied by increasing  $d$  and padding zeros to the coordinates of the sample points. Increasing  $d$  keeps a zero fitting error which allows our approach to minimize the “curviness”.

Our theoretical result should hold when each sample point is perturbed in a random direction in  $\mathbb{R}^d$  by a distance  $O(\varrho\varepsilon^3)$ , but we have not pursued the analysis as  $O(\varrho\varepsilon^3)$  is rather small. We experimented with a fair amount of noise and the estimates are satisfactory. Refer to Section 3 for details.

How does our result compare with those in the literature? In [10, 11, 15], the sample points are required to satisfy two conditions: (i) for every point  $y \in \mathcal{M}$ , the distance between  $y$  and the nearest sample point is at most the local feature size at  $y$  times  $\mu$  for some sufficiently small  $\mu \in (0, 1)$ , and (ii) for every pair of sample points  $\mathbf{p}$  and  $\mathbf{q}$ , the distance between  $\mathbf{p}$  and  $\mathbf{q}$  is at least the local feature size at  $\mathbf{p}$  times  $\delta$  for some  $\delta \in (0, \mu)$ . To estimate the tangent space at a sample point  $\mathbf{p}$ , the methods in [10, 15] use the sample points no farther from  $\mathbf{p}$  than the local feature size at  $\mathbf{p}$  times  $c\mu$  for some  $c \geq 2$ . It follows that at least  $a^m$  sample points are needed for some constant  $a > 1$  depending on  $\mathcal{M}$ . The method in [11] uses the Voronoi cell of  $\mathbf{p}$  for tangent space estimation. Using local information only, it is impossible to obtain the Voronoi cell of  $\mathbf{p}$ , and it is unclear to obtain an appropriate approximate Voronoi cell. In the worst case,  $\mathbf{p}$  can have at least  $a^m$  Voronoi neighbors for some constant  $a > 1$  depending on  $\mathcal{M}$ . The angular error bounds given in [10, 11, 15] are  $O(\mu)$ . The running times are  $O(d2^{O(m^7 \log m)})$  in [15],  $O(d2^{O(m)})$  in [10], and  $O(N^{\lceil (d+1)/2 \rceil})$  in [11], where  $N$  is the total number of sample points. In [2], the sampling is required to satisfy the condition that for every point  $y \in \mathcal{M}$ , the distance between  $y$  and the nearest sample point is at most  $\mu$  for some sufficiently small  $\mu \in (0, 1)$ . The tangent space at a sample point  $\mathbf{p}$  is estimated using the sample points within a distance  $r$  from  $\mathbf{p}$ , where  $r$  can be any value in  $[10\mu, 1/2)$ . Thus, at least  $a^m$  sample points are needed for some constant  $a > 1$  depending on  $\mathcal{M}$ . The angular error is  $O(r/\varrho)$ . The running time is  $O(dn^{O(m^6 \log m)})$ , where  $n$  is the number of sample points in the neighborhood. The work by Little et al. [26] is a multiscale analysis of the local covariance matrix. Noise is allowed and only roughly  $O(m \log m)$  points in a local neighborhood are required for computation.

In our case, although  $n$  can be as small as  $\binom{m+1}{2} + 1$  for the algorithm to be applied, the probability bound  $1 - O(n^{-1/3})$  is only meaningful for larger values of  $n$  because the hidden constant in the probability bound is a polynomial in  $m$ . Nevertheless, a polynomial in  $m$  is asymptotically smaller than  $a^m$  for any constant  $a > 1$ . This makes our neighborhood radius smaller than those in [2, 10, 11, 15] for large  $m$ , but we require more sample points than the approach in [26]. The angular error bounds in [2, 10, 11, 15] are  $O(r/\varrho)$ . (Note that  $r/\varrho < 1$ .) Roughly speaking, the angular error bound in [26] is linear in  $r$ , but the bound has a sophisticated expression and the reader is referred to [26] for details. Our angular error bound is  $O(\varepsilon^2) = O((r/\varrho)^2)$ . The hidden constant in our angular error bound depends on  $\mathcal{M}$  and a polynomial in  $m$ .

The probability bound  $1 - O(n^{-1/3})$  appears in many places in our analysis, where it is also implicitly assumed that  $n$  is greater than or equal to some appropriate polynomial in  $m$ . In practice, we suggest setting  $n \geq \binom{m+1}{2} + m$  because our approach is based on locally fitting a quadratic function, and the minimum number of variables in such a quadratic function is  $\binom{m+1}{2} + m$  when  $d = m + 1$ . In our experiments (Section 3), setting  $n = \binom{m+1}{2} + m + 30$  gives

good results in both the noiseless and noisy cases.

## 2 Problem formulation, algorithm and overview

### 2.1 Modeling

We discuss in this section how to model the local neighborhood of the origin using some implicit functions  $F_\ell$ ,  $\ell \in [m+1, d]$ , with domain  $\mathbb{R}^d$ . Recall that  $\{\mathbf{a}_p : p \in [1, n]\}$  are the  $n$  sample points nearest to the origin (which is also a sample point). The functions  $F_\ell$  are constructed so that  $F_\ell(\mathbf{a}_p) = 0$  for every  $p \in [1, n]$  and every  $\ell \in [m+1, d]$ . The goal is to obtain the compact representation of  $(F_\ell(\mathbf{a}_1) \cdots F_\ell(\mathbf{a}_n))$  in (2.3) below.

Let  $\gamma_1, \dots, \gamma_d$  be any  $d$  orthogonal coordinate axes of  $\mathbb{R}^d$  (with the same origin) such that  $\gamma_1, \dots, \gamma_m$  span  $\mathcal{T}$ .

First, we apply an (unknown) orthonormal transformation  $\psi$  so that the coordinates of each point  $\psi(\mathbf{a}_p)$  is expressed with respect to the coordinate system  $(\gamma_1, \dots, \gamma_d)$ . Recall that the  $\ell$ -th coordinate of  $\psi(\mathbf{a}_p)$ ,  $\ell \in [m+1, d]$ , is the value of the coordinate function  $f_\ell$  on the first  $m$  coordinates of  $\psi(\mathbf{a}_p)$ . Figure 1(a) shows an example of a manifold and the coordinate system before applying the transformation  $\psi$ . Figure 1(b) shows the corresponding  $f_2$  and  $f_3$  after applying  $\psi$ .

For every positive integer  $k$ , let  $D^k f_\ell|_0$  denote the  $k$ -th derivative of  $f_\ell$  at the origin, which is a map that sends  $k$  vectors from  $\mathbb{R}^m$  to a real number. The domain of  $D^k f_\ell|_0$  consists of  $k$  copies of  $\mathbb{R}^m$  spanned by  $(\gamma_1, \dots, \gamma_m)$ . When  $k = 2$ , one can view  $D^2 f_\ell|_0$  as an  $m \times m$  matrix, and then  $D^2 f_\ell|_0(\mathbf{v}, \mathbf{v})$  is equal to  $\mathbf{v}^t \cdot D^2 f_\ell|_0 \cdot \mathbf{v}$  for every vector  $\mathbf{v} \in \mathbb{R}^m$ . The matrix  $D^2 f_\ell|_0$  is known as the *Hessian matrix*.

For every vector  $\mathbf{v} \in \mathbb{R}^m$  with a small enough  $\|\mathbf{v}\|$ , the Taylor expansion of  $f_\ell(\mathbf{v})$  is

$$f_\ell(\mathbf{v}) = \frac{1}{2} D^2 f_\ell|_0(\mathbf{v}, \mathbf{v}) + \frac{1}{6} D^3 f_\ell|_0(\mathbf{v}, \mathbf{v}, \mathbf{v}) + \dots$$

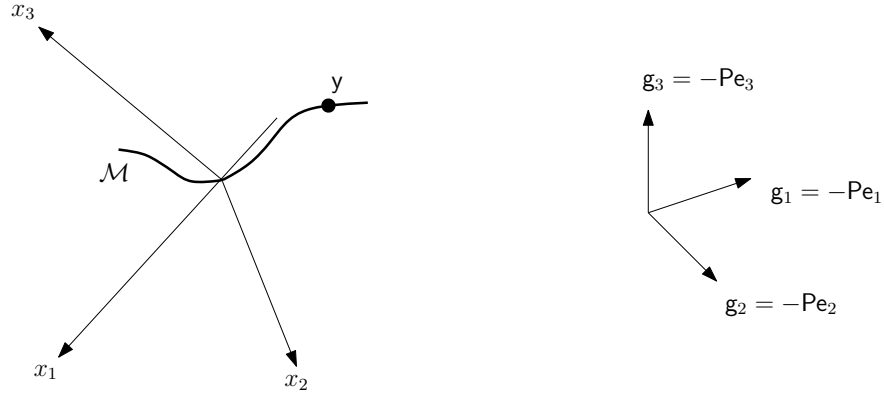
There is no constant term in the Taylor expansion above because  $\mathcal{M}$  passes through the origin. There is no linear term because  $(\gamma_1, \dots, \gamma_m)$  span  $\mathcal{T}$  and so  $Df_\ell|_0$  vanishes. We extend the domain of  $f_\ell$  from  $\mathbb{R}^m$  to  $\mathbb{R}^d$  by ignoring the last  $d-m$  coordinates of the input vector. That is, the vector  $\mathbf{v}$  can be paired with any vector  $\mathbf{w} \in \mathbb{R}^{d-m}$  to yield the following extended expansion:

$$\frac{1}{2} (\mathbf{v}^t \ \mathbf{w}^t) \begin{pmatrix} D^2 f_\ell|_0 & \mathbf{0}_{m, d-m} \\ \mathbf{0}_{d-m, m} & \mathbf{0}_{d-m, d-m} \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix} + \dots$$

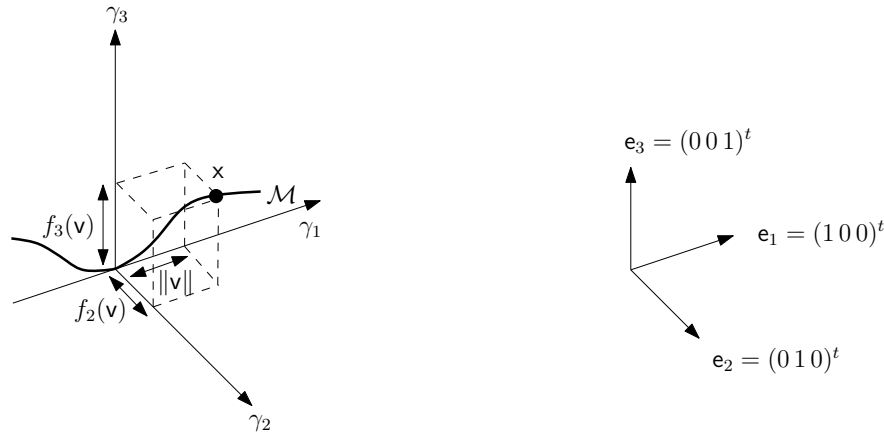
Transforming this extended expansion back to the coordinate system  $(x_1, \dots, x_d)$  gives the following function  $F_\ell : \mathbb{R}^d \rightarrow \mathbb{R}$  whose zero-set contains the origin and  $\mathbf{a}_p$  for  $p \in [1, n]$ .

$$\begin{aligned} \forall \ell \in [m+1, d], \quad F_\ell(\mathbf{y}) &= \mathbf{y}^t \mathbf{g}_\ell + \frac{1}{2} \mathbf{y}^t \mathbf{Q}_\ell \mathbf{y} + \dots, \quad \text{such that} \\ \mathbf{g}_\ell &= -\mathbf{P} (\mathbf{0}_{1, \ell-1} \ 1 \ \mathbf{0}_{1, d-\ell})^t, \\ \mathbf{Q}_\ell &= \mathbf{P} \begin{pmatrix} D^2 f_\ell|_0 & \mathbf{0}_{m, d-m} \\ \mathbf{0}_{d-m, m} & \mathbf{0}_{d-m, d-m} \end{pmatrix} \mathbf{P}^t. \end{aligned} \tag{2.1}$$

The matrix  $\mathbf{P}^t$  is the unknown  $d \times d$  orthonormal matrix that realizes the transformation  $\psi$ . The vector  $(\mathbf{0}_{1, \ell-1} \ 1 \ \mathbf{0}_{1, d-\ell})^t$  is the “vertical direction” for  $f_\ell$  at the origin with respect to the coordinate system  $(\gamma_1 \cdots \gamma_d)$ . It means that  $\mathbf{g}_\ell$  is the gradient of  $F_\ell$  and also a normal vector to  $\mathcal{M}$  at the origin.  $\mathbf{Q}_\ell$  is a  $d \times d$  matrix.



(a)



(b)

Figure 1: The manifold  $\mathcal{M}$  is a curve shown in bold in (a) and (b). The coordinate system in (a) is  $(x_1, x_2, x_3)$ , whereas  $(\gamma_1, \gamma_2, \gamma_3)$  is the coordinate system in (b). The right figure in (b) shows an orthonormal basis  $(\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3)$  of the coordinate system  $(\gamma_1, \gamma_2, \gamma_3)$ . The left figure in (a) is mapped by the orthonormal transformation  $\psi$  to the left figure in (b). The transformation  $\psi$  is realized by the orthonormal matrix  $\mathbf{P}^t$ . Therefore,  $\mathbf{x} = \mathbf{P}^t \mathbf{y}$ . Each basis vector  $\mathbf{e}_i$  with respect to the coordinate system  $(\gamma_1, \gamma_2, \gamma_3)$  is mapped to the vector  $\mathbf{g}_i = -\mathbf{P}\mathbf{e}_i$  with respect to the coordinate system  $(x_1, x_2, x_3)$ . The tangent space  $\mathcal{T}$  of  $\mathcal{M}$  at the origin is spanned by the vector  $\mathbf{g}_1$  with respect to the coordinate system  $(x_1, x_2, x_3)$  in (a), or equivalently, the vector  $\mathbf{e}_1$  with respect to the coordinate system  $(\gamma_1, \gamma_2, \gamma_3)$  in (b).

Let  $q_{\ell,ij}$  denote the  $(i,j)$  entry of  $\mathbf{Q}_\ell$ . Since  $\mathbf{Q}_\ell$  is symmetric,  $q_{\ell,ij}$  and  $q_{\ell,ji}$  are equal. Expanding the terms  $\mathbf{a}_p^t \mathbf{g}_\ell$  and  $\frac{1}{2} \mathbf{a}_p^t \mathbf{Q}_\ell \mathbf{a}_p$  in  $F_\ell(\mathbf{a}_p)$  gives:

$$\begin{aligned} F_\ell(\mathbf{a}_p) &= (a_{p1} \ a_{p2} \ \cdots \ a_{pd}) \cdot \mathbf{g}_\ell + \\ &\quad \left( \frac{1}{\sqrt{2}} a_{p1}^2 \quad a_{p1}a_{p2} \quad \cdots \quad a_{p1}a_{pd} \quad \frac{1}{\sqrt{2}} a_{p2}^2 \quad a_{p2}a_{p3} \quad \cdots \quad a_{p2}a_{pd} \quad \cdots \quad \frac{1}{\sqrt{2}} a_{pd}^2 \right) \cdot \\ &\quad \left( \frac{1}{\sqrt{2}} q_{\ell,11} \quad q_{\ell,12} \quad \cdots \quad q_{\ell,1d} \quad \frac{1}{\sqrt{2}} q_{\ell,22} \quad q_{\ell,23} \quad \cdots \quad q_{\ell,2d} \quad \cdots \quad \frac{1}{\sqrt{2}} q_{\ell,dd} \right)^t + \cdots \end{aligned}$$

This motivates us to define:

$$\begin{aligned} \mathbf{c}_\ell &\stackrel{\text{def}}{=} \left( \frac{1}{\sqrt{2}} q_{\ell,11} \quad q_{\ell,12} \quad \cdots \quad q_{\ell,1d} \quad \frac{1}{\sqrt{2}} q_{\ell,22} \quad q_{\ell,23} \quad \cdots \quad q_{\ell,2d} \quad \cdots \quad \frac{1}{\sqrt{2}} q_{\ell,dd} \right)^t \\ \mathbf{A} &\stackrel{\text{def}}{=} \begin{pmatrix} a_{11} & \cdots & a_{1d} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nd} \end{pmatrix} \\ \mathbf{B} &\stackrel{\text{def}}{=} \begin{pmatrix} \frac{1}{\sqrt{2}} a_{11}^2 & a_{11}a_{12} & \cdots & a_{11}a_{1d} & \frac{1}{\sqrt{2}} a_{12}^2 & a_{12}a_{13} & \cdots & \frac{1}{\sqrt{2}} a_{1d}^2 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{2}} a_{n1}^2 & a_{n1}a_{n2} & \cdots & a_{n1}a_{nd} & \frac{1}{\sqrt{2}} a_{n2}^2 & a_{n2}a_{n3} & \cdots & \frac{1}{\sqrt{2}} a_{nd}^2 \end{pmatrix} \end{aligned} \tag{2.2}$$

$\mathbf{A}$  is an  $n \times d$  matrix,  $\mathbf{B}$  is an  $n \times \binom{d+1}{2}$  matrix, and  $\mathbf{c}_\ell$  is a  $\binom{d+1}{2}$ -dimensional vector. They yield:

$$\begin{pmatrix} F_\ell(\mathbf{a}_1) \\ \vdots \\ F_\ell(\mathbf{a}_n) \end{pmatrix} = (\mathbf{A} \ \mathbf{B}) \cdot \begin{pmatrix} \mathbf{g}_\ell \\ \mathbf{c}_\ell \end{pmatrix} + \cdots \tag{2.3}$$

The coefficient  $1/\sqrt{2}$  of  $a_{pi}^2$  in the definition of  $\mathbf{B}$  is needed so that the eigenvalues of  $\mathbf{B}^t \mathbf{B}$  are independent of rotations in  $\mathbb{R}^d$  that keep the origin fixed. We will establish this fact in Lemma 6.2. The vectors  $\mathbf{g}_\ell$  and  $\mathbf{c}_\ell$  are unknowns in (2.3), and  $\mathbf{c}_\ell$  is a linearization of the matrix  $\mathbf{Q}_\ell$ . Since  $F_\ell^{-1}(0)$  contains  $\mathbf{a}_p$  for  $p \in [1, n]$ , the left hand side of (2.3) is a zero vector.

To approximate the  $F_\ell$ 's using quadratic functions, the first attempt is to retain just  $(\mathbf{A} \ \mathbf{B}) \cdot (\mathbf{g}_\ell^t \ \mathbf{c}_\ell^t)^t$  in the right hand side of (2.3) because this keeps only the linear and quadratic terms. The subsequent analysis in Section 7 demands  $\|\mathbf{B}^\dagger\|$  to be comparable to the reciprocal of the  $(\binom{m+1}{2} + 1)$ -th largest singular value of  $\mathbf{B}$ . Unfortunately,  $\|\mathbf{B}^\dagger\|$  is determined by the possibly much larger reciprocal of the smallest singular value of  $\mathbf{B}$ . As a result, we modify  $\mathbf{B}$  by changing its  $n - \binom{m+1}{2}$  smallest singular values as follows. Recall that  $\binom{m+1}{2} + 1 \leq n \leq \binom{d+1}{2}$  by assumption. Define:

$$\begin{aligned} m_0 &\stackrel{\text{def}}{=} \binom{m+1}{2} \\ \mathbf{L}\mathbf{A}\mathbf{R}^t &\stackrel{\text{def}}{=} \text{thin SVD of } \mathbf{B}, \text{ where } \mathbf{\Lambda} = \text{diag}_n(\lambda_1, \lambda_2, \dots, \lambda_n) \text{ and } \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \geq 0 \\ \widehat{\mathbf{\Lambda}} &\stackrel{\text{def}}{=} \text{diag}_n(\lambda_1, \lambda_2, \dots, \lambda_{m_0}, \underbrace{\lambda_{m_0+1}, \dots, \lambda_{m_0+1}}_{n-m_0 \text{ copies}}) \\ \widehat{\mathbf{B}} &\stackrel{\text{def}}{=} \widehat{\mathbf{L}}\widehat{\mathbf{A}}\mathbf{R}^t \end{aligned}$$

$\mathbf{L}$  is an  $n \times n$  matrix,  $\mathbf{\Lambda}$  and  $\widehat{\mathbf{\Lambda}}$  are  $n \times n$  diagonal matrices, and  $\mathbf{R}$  is a  $\binom{d+1}{2} \times n$  matrix.  $\widehat{\mathbf{B}}$  is the replacement of  $\mathbf{B}$ .



Lemma 2.1 below shows that the case of  $\lambda_{m_0+1} = 0$  can be dealt with separately in the algorithm. (The proof of Lemma 2.1 is given in Section 8.) Thus, we can assume that  $\lambda_{m_0+1}$  is positive. For ease of presentation, we also assume that  $\lambda_n > 0$ ; otherwise, for each  $\lambda_i = 0$  (such an  $i$  can range from  $m_0 + 2$  to  $n$ ), we set the corresponding diagonal entry of  $\widehat{\Lambda}$  to zero. The proof of Lemma 7.5 requires that if a diagonal entry of  $\Lambda$  is zero, the corresponding entry in  $\widehat{\Lambda}$  is also zero.

**Lemma 2.1** *There exists a constant  $c$  that is a polynomial in  $m$  such that if  $n \geq c$ ,  $\varepsilon$  is sufficiently small, and  $\lambda_{m_0+1} = 0$ , then with probability  $1 - O(n^{-1/3})$ ,  $\mathcal{T}$  is equal to the space spanned by the eigenvectors corresponding to the  $m$  largest eigenvalues of  $A^t A$ .*

By keeping only  $(A \ B) \cdot (\mathbf{g}_\ell^t \ \mathbf{c}_\ell^t)^t$  in the right hand side of (2.3) and replacing  $B$  by  $\widehat{B}$ , we allude to some nonlinear functions  $\widehat{F}_\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\ell \in [m + 1, d]$ , that approximate the  $F_\ell$ 's and satisfy the following system:

$$\forall \ell \in [m + 1, d], \quad \begin{pmatrix} \widehat{F}_\ell(\mathbf{a}_1) \\ \vdots \\ \widehat{F}_\ell(\mathbf{a}_n) \end{pmatrix} = (A \ \widehat{B}) \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell \end{pmatrix}, \quad (2.4)$$

where  $\widehat{\mathbf{g}}_\ell$  and  $\widehat{\mathbf{c}}_\ell$  are the new unknowns. The vector  $\widehat{\mathbf{g}}_\ell$  is the unknown gradient of  $\widehat{F}_\ell$ . Since the unknown gradients are supposed to span the approximate normal space at the origin, we require them to satisfy the following constraints:

$$\begin{aligned} \forall \ell \in [m + 1, d], \quad & \|\widehat{\mathbf{g}}_\ell\| = 1 \\ \forall \ell_1 \neq \ell_2, \quad & \widehat{\mathbf{g}}_{\ell_1} \perp \widehat{\mathbf{g}}_{\ell_2} \end{aligned} \quad (2.5)$$

## 2.2 Algorithm and overview

Our algorithm solves an eigenvalue problem derived from  $\widehat{B}$ . Define the following matrix:

$$H = (\widehat{B}^\dagger A)^t.$$

We find the eigenvectors corresponding to the  $m$  largest eigenvalues of the following matrix:

$$HH^t = A^t (\widehat{B}^\dagger)^t \cdot \widehat{B}^\dagger A = A^t L \widehat{\Lambda}^\dagger \widehat{\Lambda}^\dagger L^t A.$$

These  $m$  eigenvectors span the approximate tangent space. We compute the full SVD of the  $n \times d$  matrix  $\widehat{\Lambda}^\dagger L^t A$  to obtain the eigenvectors of  $HH^t$ . The pseudocode is given below.

TANGENT(A)

1. Compute the thin SVD  $LAR^t$  of  $B$ .
2. If  $\lambda_{m_0+1} = 0$ , then return the eigenvectors corresponding to the  $m$  largest eigenvalues of  $A^t A$  as an orthonormal basis of the estimated tangent space.
3. Compute  $\widehat{\Lambda}^\dagger$  and the full SVD  $CDE^t$  of  $\widehat{\Lambda}^\dagger L^t A$ . Assume that the diagonal entries of  $D$  are in descending order and that  $E = (\widehat{\mathbf{g}}_1, \dots, \widehat{\mathbf{g}}_d)$ , where  $\widehat{\mathbf{g}}_\ell$  corresponds to the  $\ell$ th largest diagonal entry in  $D$ .
4. Return  $(\widehat{\mathbf{g}}_1, \dots, \widehat{\mathbf{g}}_m)$  as an orthonormal basis of the approximate tangent space.

The running time is  $O(T_{\text{tsvd}}(n, \binom{d+1}{2}) + T_{\text{fsvd}}(n, d) + d^2 n)$ , which is  $O(d^2 n^2)$  in practice [7, 16, 17].

Section 3 describes our experiments that demonstrate the accuracy of our tangent estimation in both noiseless and noisy cases.

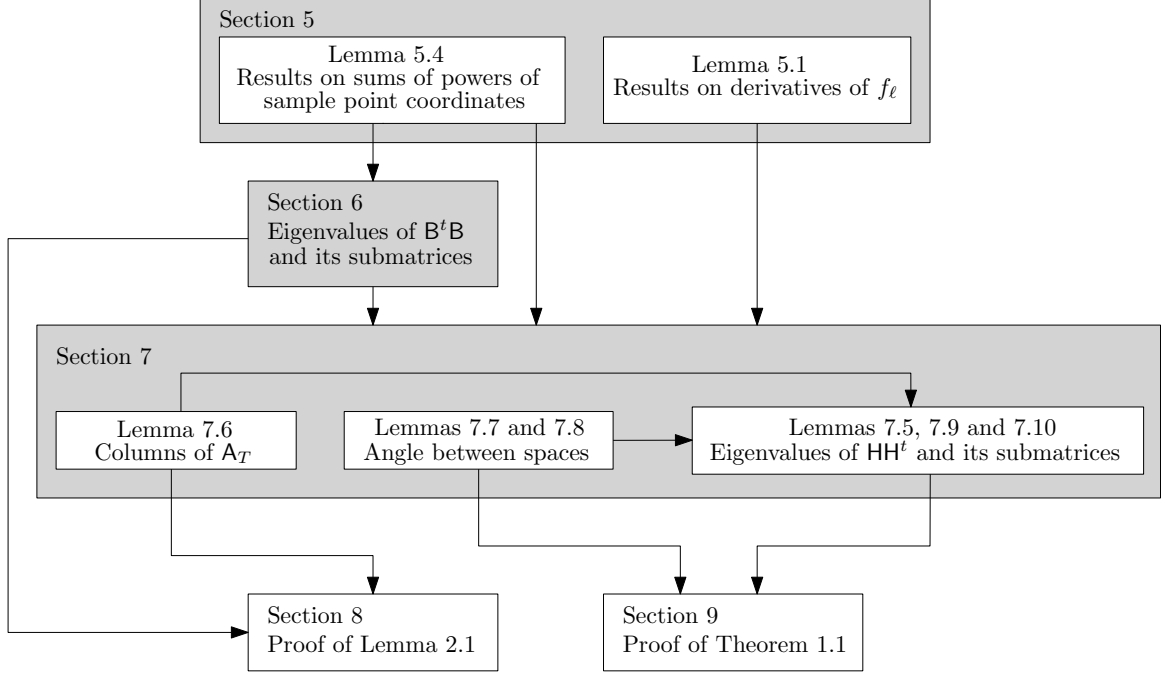


Figure 2: Proof overview.

Section 4 turns (2.4) and (2.5) into an eigenvalue problem. The solution minimizes the quantity  $\frac{1}{n} \sum_{\ell=m+1}^d \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$ , that is, the sum of squared interpolation error with the “curviness” penalized by the term  $\sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$ . Tuning the parameter  $\alpha$  is time-consuming though. We push  $\alpha$  to zero in the limit to derive the eigenvalue problem for  $\mathbf{H}\mathbf{H}^t$ .

Sections 5–9 contain the analysis that leads to an  $O(\varepsilon^2)$  bound on the angular error. Figure 2 shows the dependence of results in different sections. Assume that the coordinate axes  $x_1, \dots, x_m$  span the tangent space of  $\mathcal{M}$  at the origin. So the coordinate axes  $x_{m+1}, \dots, x_d$  span the normal space at the origin.

Section 5 presents several results on the derivatives of the coordinate functions and on the sums of powers of sample point coordinates. These results may be of independent interest. Lemma 5.1(i) shows that  $\|(\mathbf{D}^2 f_{m+1}|_0(\mathbf{u}, \mathbf{u}) \cdots \mathbf{D}^2 f_d|_0(\mathbf{u}, \mathbf{u}))\| = O(1/\varrho)$  for any unit vector  $\mathbf{u} \in \mathbb{R}^m$ , independent of the dimension and other factors. Similarly, Lemma 5.1(ii) shows that  $(f_{m+1}(\mathbf{v}), \dots, f_d(\mathbf{v}))$  is approximated well by  $(\frac{1}{2}\mathbf{D}^2 f_{m+1}|_0(\mathbf{v}, \mathbf{v}), \dots, \frac{1}{2}\mathbf{D}^2 f_d|_0(\mathbf{v}, \mathbf{v}))$  for any  $\mathbf{v} \in \mathbb{R}^m$  such that  $\|\mathbf{v}\| \leq \varrho\varepsilon$ , that is,  $\|(f_{m+1}(\mathbf{v}) - \frac{1}{2}\mathbf{D}^2 f_{m+1}|_0(\mathbf{v}, \mathbf{v}), \dots, f_d(\mathbf{v}) - \frac{1}{2}\mathbf{D}^2 f_d|_0(\mathbf{v}, \mathbf{v}))\| = O(\varrho\varepsilon^3)$ . Lemma 5.4 gives concentration bounds for  $\left| \sum_{p=1}^n a_{pi} a_{pj} a_{pk} a_{pl} \right|$ ,  $\left| \sum_{p=1}^n a_{pi} a_{pj} a_{pk} \right|$ , and  $\left| \sum_{p=1}^n a_{pi} a_{pj} \right|$  for  $i, j, k, l \in [1, m]$ .

Since the approximate tangent space is spanned by some eigenvectors of  $\mathbf{H}\mathbf{H}^t$ , we need to analyze the eigenvalues of  $\mathbf{H}\mathbf{H}^t = \mathbf{A}^t(\widehat{\mathbf{B}}^\dagger)^t \cdot \widehat{\mathbf{B}}^\dagger \mathbf{A}$ , which requires us to bound the eigenvalues of  $\widehat{\mathbf{B}}^\dagger$ . These results are presented in Section 6. Rearrange the columns of  $\mathbf{B}$  so that  $\mathbf{B} = (\mathbf{B}_{TT} \ \mathbf{B}_{TN} \ \mathbf{B}_{NN})$ , where  $\mathbf{B}_{TT}$  consists of columns in  $a_{pi} a_{pj}$  for possibly non-distinct  $i, j \in [1, m]$ ,  $\mathbf{B}_{TN}$  consists of columns in  $a_{pi} a_{pj}$  for  $i \in [1, m]$  and  $j \in [m+1, d]$ , and  $\mathbf{B}_{NN}$  consists of columns

in  $a_{pi}a_{pj}$  for possibly non-distinct  $i, j \in [m+1, d]$ . We divide  $\mathbf{B}^t\mathbf{B}$  into blocks as follows.

$$\mathbf{B}^t\mathbf{B} = \begin{pmatrix} \mathbf{B}_{TT}^t\mathbf{B}_{TT} & \mathbf{B}_{TT}^t\mathbf{B}_{TN} & \mathbf{B}_{TT}^t\mathbf{B}_{NN} \\ \mathbf{B}_{TN}^t\mathbf{B}_{TT} & \mathbf{B}_{TN}^t\mathbf{B}_{TN} & \mathbf{B}_{TN}^t\mathbf{B}_{NN} \\ \mathbf{B}_{NN}^t\mathbf{B}_{TT} & \mathbf{B}_{NN}^t\mathbf{B}_{TN} & \mathbf{B}_{NN}^t\mathbf{B}_{NN} \end{pmatrix}$$

We apply the concentration bounds in Lemma 5.4. Then, the Gershgorin Circle Theorem [13, 17] says that the eigenvalues of  $\mathbf{B}_{TT}^t\mathbf{B}_{TT}$  are dominated by the diagonal entries which are  $\Theta(n\rho^4\varepsilon^4)$  (Lemma 6.3). Using known bounds on  $\sum_{i=1}^m a_{pi}^2$  and  $\sum_{i=m+1}^d a_{pi}^2$ , one can easily show that  $\|\mathbf{B}_{TN}\| = O(\sqrt{n}\rho^2\varepsilon^3)$  and  $\|\mathbf{B}_{NN}\| = O(\sqrt{n}\rho^2\varepsilon^4)$  (Lemma 6.4). Then, applying the Gershgorin Circle Theorem to the division of  $\mathbf{B}^t\mathbf{B}$  above shows that the  $m_0$  largest eigenvalues of  $\mathbf{B}^t\mathbf{B}$  are dominated by those of  $\mathbf{B}_{TT}^t\mathbf{B}_{TT}$  and hence are  $\Theta(n\rho^4\varepsilon^4)$  (Lemma 6.5). A finer analysis then shows that the  $(m_0+1)$ -th largest eigenvalue  $\lambda_{m_0+1}^2$  of  $\mathbf{B}^t\mathbf{B}$  is  $O(n\rho^4\varepsilon^6)$  (Lemma 6.6). Note that  $\|\widehat{\mathbf{B}}^\dagger\| = 1/\lambda_{m_0+1}$ .

We write  $\mathbf{A} = (\mathbf{A}_T \mathbf{A}_N)$ , where  $\mathbf{A}_T$  consists of the columns  $(a_{1i}, \dots, a_{mi})^t$  for  $i \in [1, m]$  and  $\mathbf{A}_N$  consists of the columns  $(a_{1i}, \dots, a_{ni})^t$  for  $i \in [m+1, d]$ . Since  $\mathbf{H} = (\widehat{\mathbf{B}}^\dagger\mathbf{A})^t$ , we can write

$$\mathbf{H}_T = (\widehat{\mathbf{B}}^\dagger\mathbf{A}_T)^t, \quad \mathbf{H}_N = (\widehat{\mathbf{B}}^\dagger\mathbf{A}_N)^t, \quad \mathbf{H} = \begin{pmatrix} \mathbf{H}_T \\ \mathbf{H}_N \end{pmatrix}, \quad \mathbf{H}\mathbf{H}^t = \begin{pmatrix} \mathbf{H}_T\mathbf{H}_T^t & \mathbf{H}_T\mathbf{H}_N^t \\ \mathbf{H}_N\mathbf{H}_T^t & \mathbf{H}_N\mathbf{H}_N^t \end{pmatrix}.$$

Note that  $\mathbf{H}_T$  is an  $m \times \binom{d+1}{2}$  submatrix and  $\mathbf{H}_N$  is a  $(d-m) \times \binom{d+1}{2}$  submatrix. In Section 7, we prove that  $\|\mathbf{H}_N\| = O(\sqrt{n}\rho\varepsilon^3/\lambda_{m_0+1})$  and every singular value of  $\mathbf{H}_T$  is  $\Theta(\sqrt{n}\rho\varepsilon/\lambda_{m_0+1})$ . Then, the Gershgorin Circle Theorem implies that the  $m$  largest eigenvalues of  $\mathbf{H}\mathbf{H}^t$  are dominated by those of  $\mathbf{H}_T\mathbf{H}_T^t$  and hence they are  $\Theta(n\rho^2\varepsilon^2/\lambda_{m_0+1}^2)$ , and the  $d-m$  smallest eigenvalues of  $\mathbf{H}\mathbf{H}^t$  are at most  $\|\mathbf{H}_N\|^2 + \|\mathbf{H}_T\|\|\mathbf{H}_N\| = O(n\rho^2\varepsilon^4/\lambda_{m_0+1}^2)$  (Lemma 7.10).

In the analysis of  $\mathbf{H}_N^t = \widehat{\mathbf{B}}^\dagger\mathbf{A}_N$  (Lemma 7.5), our definition of  $\widehat{\mathbf{B}}$  gives  $\|\widehat{\mathbf{B}}^\dagger\mathbf{A}_N\| \leq \|\widehat{\mathbf{B}}^\dagger\mathbf{A}_N\|$ . Notice that  $\mathbf{B}(\widehat{\mathbf{B}}^\dagger\mathbf{A}_N) = \mathbf{A}_N$ . By the Taylor expansion, we can approximate each entry  $a_{pl}$  in  $\mathbf{A}_N$  by  $\frac{1}{2}(a_{p1} \cdots a_{pm}) \cdot \mathbf{D}^2 f_\ell|_0 \cdot (a_{p1} \cdots a_{pm})^t$ . The definition of  $\mathbf{B}$  allows us to write these Taylor expansions as  $\mathbf{B}\mathbf{Z} \approx \mathbf{A}_N$ , where  $\mathbf{Z}$  is a  $\binom{d+1}{2} \times (d-m)$  matrix such that the  $(\ell-m)$ -th column is

$$\left( \frac{1}{\sqrt{2}}q_{\ell,11} \quad q_{\ell,12} \quad \cdots \quad q_{\ell,1d} \quad \frac{1}{\sqrt{2}}q_{\ell,22} \quad q_{\ell,23} \quad \cdots \quad q_{\ell,2d} \quad \cdots \quad \frac{1}{\sqrt{2}}q_{\ell,dd} \right)^t,$$

where  $q_{\ell,ij}$  is the  $(i, j)$  entry of the Hessian matrix  $\mathbf{D}^2 f_\ell|_0$  if  $i, j \in [1, m]$ , and  $q_{\ell,ij} = 0$  otherwise. Then, Lemma 5.1(i) allows us to conclude that  $\|\mathbf{Z}\| = O(1/\rho)$ , and by the property of pseudoinverse,  $\|\widehat{\mathbf{B}}^\dagger\mathbf{A}_N\| \lesssim \|\mathbf{Z}\| = O(1/\rho) = O(\sqrt{n}\rho\varepsilon^3/\lambda_{m_0+1})$  as  $\lambda_{m_0+1} = O(\sqrt{n}\rho^2\varepsilon^3)$  by the result in Section 6.

The analysis of  $\mathbf{H}_T$  makes use of the thin SVD  $\widehat{\mathbf{L}}\widehat{\mathbf{A}}\mathbf{R}^t$  of  $\widehat{\mathbf{B}}$  (Lemma 7.9). Recall that the  $m_0$  largest singular values of  $\widehat{\mathbf{B}}$  are those of  $\mathbf{B}$ . The other  $n-m_0$  singular values are equal to  $\lambda_{m_0+1}$ . If we group the largest  $m_0$  singular values in an  $m_0 \times m_0$  diagonal submatrix  $\widehat{\Lambda}_0$ , then we can write  $\widehat{\Lambda}$ ,  $\widehat{\mathbf{L}}$ , and  $\mathbf{R}$  as follows.

$$\widehat{\Lambda} = \begin{pmatrix} \widehat{\Lambda}_0 & 0_{m_0, n-m_0} \\ 0_{n-m_0, m_0} & \lambda_{m_0+1} \mathbf{I}_{n-m_0} \end{pmatrix}, \quad \widehat{\mathbf{L}} = \left( \underbrace{\mathbf{L}_0}_{m_0 \text{ columns}} \quad \underbrace{\mathbf{L}_1}_{n-m_0 \text{ columns}} \right), \quad \mathbf{R} = \left( \underbrace{\mathbf{R}_0}_{m_0 \text{ columns}} \quad \underbrace{\mathbf{R}_1}_{n-m_0 \text{ columns}} \right).$$

Combining the above with the relation  $\mathbf{H}_T^t = \widehat{\mathbf{B}}^\dagger\mathbf{A}_T$ , we obtain

$$\mathbf{H}_T\mathbf{H}_T^t = \mathbf{A}_T^t\mathbf{L}_0 \left( \widehat{\Lambda}_0^\dagger \right)^2 \mathbf{L}_0^t\mathbf{A}_T + \frac{1}{\lambda_{m_0+1}^2} \mathbf{A}_T^t\mathbf{L}_1\mathbf{L}_1^t\mathbf{A}_T.$$

Thus, the maximum eigenvalue of  $\mathbf{H}_T\mathbf{H}_T^t$  is at most  $\|\mathbf{A}_T^t\mathbf{L}_0 \left( \widehat{\Lambda}_0^\dagger \right)^2 \mathbf{L}_0^t\mathbf{A}_T\| + \frac{1}{\lambda_{m_0+1}^2} \|\mathbf{A}_T^t\mathbf{L}_1\mathbf{L}_1^t\mathbf{A}_T\|$ , which can be verified to be  $O(n\rho^2\varepsilon^2/\lambda_{m_0+1}^2)$  using the facts that  $\|\mathbf{L}_0\| = \|\mathbf{L}_1\| = 1$ ,  $\|\mathbf{A}_T\| =$

$O(\sqrt{n}\rho\varepsilon)$  as  $\|\mathbf{a}_p\| = O(\rho\varepsilon)$ ,  $\|\widehat{\Lambda}_0^\dagger\| = \Theta(1/(\sqrt{n}\rho^2\varepsilon^2))$ , and  $\lambda_{m_0+1} = O(\sqrt{n}\rho^2\varepsilon^3)$ . The minimum eigenvalue of  $\mathbf{H}_T\mathbf{H}_T^t$  is at least the minimum eigenvalue of  $\frac{1}{\lambda_{m_0+1}^2}\mathbf{A}_T^t\mathbf{L}_1\mathbf{L}_1^t\mathbf{A}_T$ . Using the concentration bounds in Lemma 5.4, we can show that every column vector  $(a_{1i}, \dots, a_{ni})^t$  in  $\mathbf{A}_T$  has a 2-norm of  $\Theta(\sqrt{n}\rho\varepsilon)$ , and that the angle between two distinct column vectors  $(a_{1i}, \dots, a_{ni})^t$  and  $(a_{1j}, \dots, a_{nj})^t$  in  $\mathbf{A}_T$  is large (Lemma 7.6). Then, we show that every column vector in  $\mathbf{A}_T$  makes a small angle with the column space of  $\mathbf{L}_1$  (Lemmas 7.7–7.8). It follows that  $\mathbf{A}_T^t\mathbf{L}_1\mathbf{L}_1^t\mathbf{A}_T \approx \|\mathbf{A}_T\|^2 = \sqrt{m} \cdot \Omega(n\rho^2\varepsilon^2)$ , and so the minimum eigenvalue of  $\mathbf{H}_T\mathbf{H}_T^t$  is at least  $\frac{1}{\lambda_{m_0+1}^2}\mathbf{A}_T^t\mathbf{L}_1\mathbf{L}_1^t\mathbf{A}_T = \Omega(n\rho^2\varepsilon^2/\lambda_{m_0+1}^2)$  as desired.

We analyze the angular error (Theorem 1.1) in Section 9 using the bounds on the singular values of  $\mathbf{H}_T$  and  $\|\mathbf{H}_N\|$ . We take an arbitrary unit eigenvector  $\mathbf{e}$  corresponding to any of the  $m$  largest eigenvalues of  $\mathbf{H}\mathbf{H}^t$ , say  $\sigma$ . Let  $\mathbf{v}$  be the vector consisting of the first  $m$  coordinates of  $\mathbf{e}$ . Let  $\mathbf{w}$  be the vector consisting of the other  $d - m$  coordinates of  $\mathbf{e}$ . We check the angle that  $\mathbf{e}$  makes the true tangent space spanned by the coordinate axes  $x_1, \dots, x_m$ . This is done by examining the equation:

$$\mathbf{H}\mathbf{H}^t\mathbf{e} = \mathbf{H}\mathbf{H}^t \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix} = \begin{pmatrix} \mathbf{H}_T\mathbf{H}_T^t & \mathbf{H}_T\mathbf{H}_N^t \\ \mathbf{H}_N\mathbf{H}_T^t & \mathbf{H}_N\mathbf{H}_N^t \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix} = \sigma \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix}.$$

Then,  $\mathbf{w} = (\sigma\mathbf{I}_{d-m} - \mathbf{H}_N\mathbf{H}_N^t)^{-1}\mathbf{H}_N\mathbf{H}_T^t\mathbf{v}$ . Therefore, the angle between  $\mathbf{e}$  and the true tangent space is  $\arctan(\|\mathbf{w}\|/\|\mathbf{v}\|) \leq \|(\sigma\mathbf{I}_{d-m} - \mathbf{H}_N\mathbf{H}_N^t)^{-1}\| \cdot \|\mathbf{H}_N\| \cdot \|\mathbf{H}_T^t\|$ . Since  $\sigma = \Theta(\|\mathbf{H}_T\|^2) = \Theta(n\rho^2\varepsilon^2/\lambda_{m_0+1}^2)$  and  $\|\mathbf{H}_N\| = O(\sqrt{n}\rho\varepsilon^3/\lambda_{m_0+1})$  by the results in Section 7, we conclude that  $\arctan(\|\mathbf{w}\|/\|\mathbf{v}\|) = O(\varepsilon^2)$ . It follows that the space spanned by the eigenvectors corresponding to the  $m$  largest eigenvalues of  $\mathbf{H}\mathbf{H}^t$  makes an  $O(\varepsilon^2)$  angle with the true tangent space, completing the analysis of the angular error.

### 3 Experimental results

We carried out some experiments to estimate tangent spaces using sample points drawn from different manifolds, including spheres, manifolds with saddles, and sinusoidal curves in high dimensions. Since our algorithm works locally, sample points are only needed in a local neighborhood. We do not know how to locally sample a manifold uniformly or according to a Poisson distribution in general, so some adhoc heuristics are used for each class of manifolds. We will describe the sampling heuristics in each case. For each manifold tested, we fix a point  $\mathbf{p}$  in the manifold and generate point samples in its neighborhood. Since we know the true tangent space, we can compute the angular error for each trial, which allows us to report the mean angular error over all trials. We conducted 25 trials for every manifold, every value of  $n$ , and every neighborhood radius. The smallest value of  $n$  is  $\binom{m+1}{2} + m + 30$  and we increase  $n$  by adding multiples of 20.

#### 3.1 Sphere

We tried unit spheres  $\mathbb{S}^m$  for  $3 \leq m \leq 9$  in  $\mathbb{R}^{200}$ . We choose  $\mathbf{p}$  to be the north pole and put an  $m$ -ball  $D$  of radius  $r$  tangent to  $\mathbb{S}^m$  at  $\mathbf{p}$ . We generate  $n$  sample points in  $D$  uniformly at random and then project them towards the origin onto  $\mathbb{S}^m$ .

The projection of  $D$  towards the origin onto  $\mathbb{S}^m$  is the spherical cap that sample points may occupy. Projecting this spherical cap orthogonally onto  $D$  gives a concentric, smaller  $m$ -ball  $D'$ . We can view  $D'$  as part of the domain of the coordinate functions  $f_{m+1}, \dots, f_d$ , which map a point in  $D'$  to a point on  $\mathbb{S}^m$ . Therefore,  $\max_{\mathbf{z} \in D'} \|f_{m+1}(\mathbf{z}), \dots, f_d(\mathbf{z})\|$  is the maximum distance between a point in  $D'$  and  $\mathbb{S}^m$ . By elementary trigonometry, this maximum distance

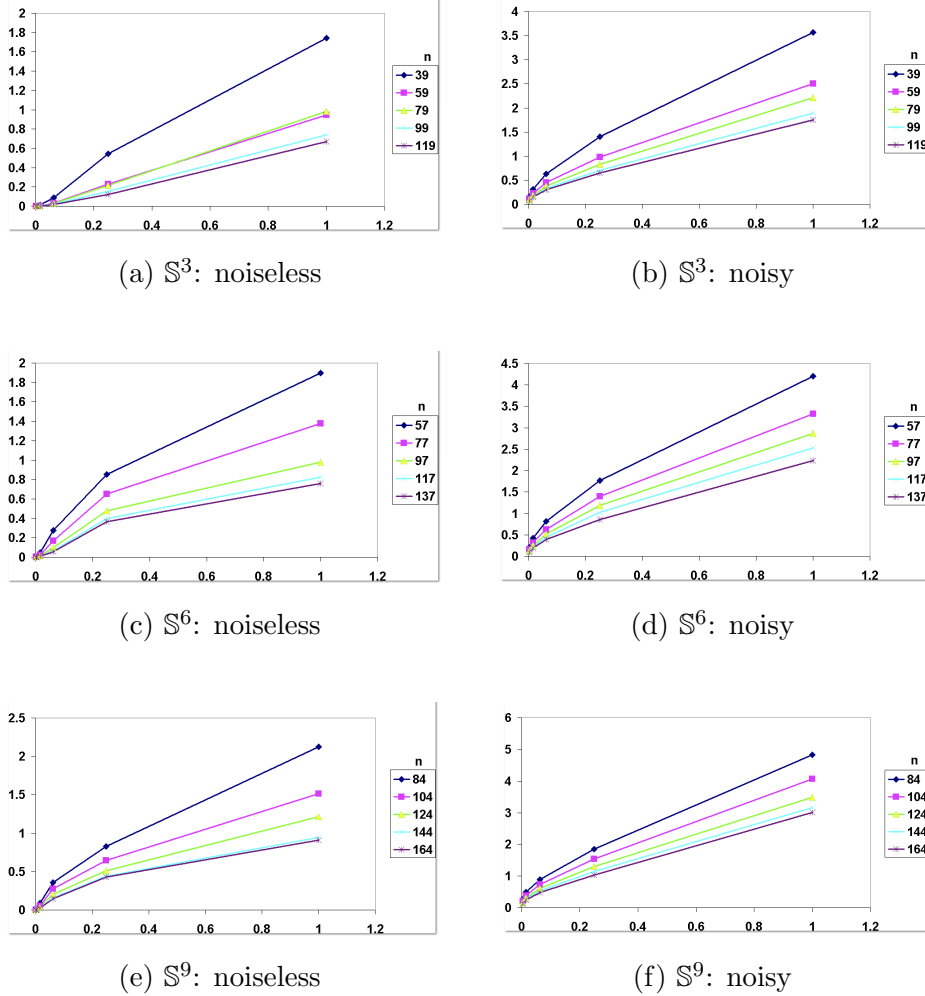


Figure 3: Plots of mean angular error for  $\mathbb{S}^3$ ,  $\mathbb{S}^6$  and  $\mathbb{S}^9$ . The vertical axes measure the mean angular errors. The horizontal axes measure  $r^2$ . The maximum noise level is at least 16%.

is  $r^2/(1 + r^2 + \sqrt{1 + r^2})$ . We vary  $r$  between 0 and 1, so  $r^2/(1 + r^2 + \sqrt{1 + r^2})$  lies between  $r^2/(2 + \sqrt{2})$  and  $r^2/2$ . Although the neighborhood radius is not exactly  $r$ , it is roughly  $cr$  for some constant  $c$ .

We plot the mean angular error against  $r^2$  for  $\mathbb{S}^3$  in Figure 3(a). A mean angular error of  $1.8^\circ$  can be achieved with  $n = 40$  even when  $r$  is as large as 1. We experimented with noisy data by adding random noise. We perturb each sample point by a random displacement chosen from  $[0, 0.08r^2]$  in a random direction in  $\mathbb{R}^d$ . Therefore, the maximum noise level is at least 16% of  $\max_{\mathbf{z} \in D'} \|f_{m+1}(\mathbf{z}), \dots, f_d(\mathbf{z})\|$ . We plot the mean angular error against  $r^2$  for  $\mathbb{S}^3$  in the noisy case in Figure 3(b). A mean angular error of roughly  $3.5^\circ$  can still be achieved with  $n = 40$  even when  $r$  is as large as 1. Both plots in the noiseless and noisy cases demonstrate that the angular error is roughly proportional to  $r^2$ . Figures 3(c)–(f) show the plots of the mean angular errors for  $\mathbb{S}^6$  and  $\mathbb{S}^9$  in the noiseless and noisy cases.

### 3.2 Manifold with saddles

Let  $S^{m-1}$  be the  $(m - 1)$ -sphere centered at the origin with radius 4 in the subspace spanned by the  $x_1, \dots, x_m$  axes. For each point  $\mathbf{q} \in S^{m-1}$ , construct the circle centered at  $\mathbf{q}$  with radius 2 embedded in the plane spanned by the vector  $\mathbf{q}$  and the  $x_{m+1}$ -axis. The union of all such

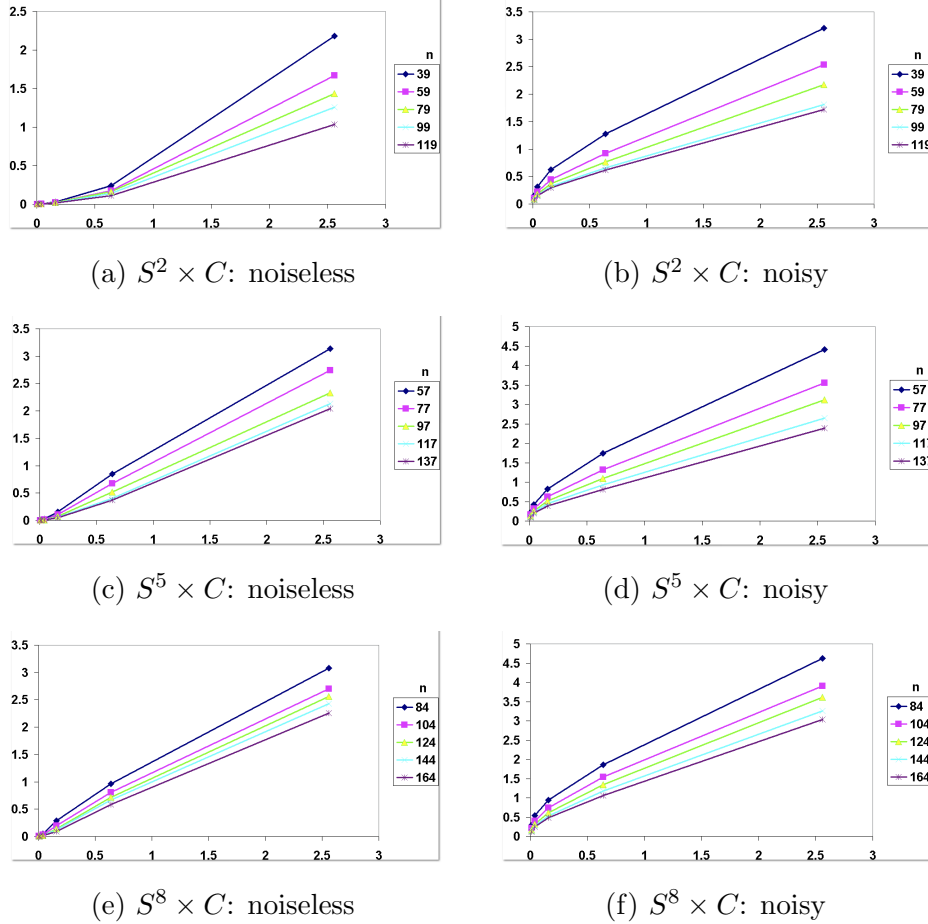


Figure 4: Plots of mean angular error for  $S^2 \times C$ ,  $S^5 \times C$  and  $S^8 \times C$ . The vertical axes measure the mean angular errors. The horizontal axes measure  $r^2$ . The maximum noise level is at least 16%.

circles is an  $m$ -dimensional manifold, which we denote by  $S^{m-1} \times C$ . We choose the point  $\mathbf{p}$  on the  $x_1$ -axis at distance 2 from the origin, which is a saddle of  $S^{m-1} \times C$ . We sample  $n$  points in an  $m$ -ball  $D$  of radius  $r$  and tangent at  $\mathbf{p}$ . Then, we lift the points orthogonally away from  $D$  and onto  $S^{m-1} \times C$ .

Since we lift sampled points in  $D$  orthogonally onto  $S^{m-1} \times C$ ,  $D$  is part of the domain of the coordinate functions  $f_{m+1}, \dots, f_d$ . The maximum distance between  $D$  and  $S^{m-1} \times C$  is  $\max_{\mathbf{z} \in D} \|f_{m+1}(\mathbf{z}), \dots, f_d(\mathbf{z})\|$ , which by elementary trigonometry is  $r^2/(2 + \sqrt{4 - r^2})$ . We vary  $r$  between 0 and 1.6, so  $r^2/(2 + \sqrt{4 - r^2})$  lies between  $0.25r^2$  and  $0.3125r^2$ . Although the neighborhood radius is not exactly  $r$ , it is roughly  $cr$  for some constant  $c$ .

We plot the mean angular error against  $r^2$  for  $S^2 \times C$  in Figure 4(a). Random noise is added in the same way as before. We perturb each sample point by a random displacement chosen from  $[0, 0.05r^2]$  in a random direction in  $\mathbb{R}^d$ . Therefore, the maximum noise level is at least 16% of  $\max_{\mathbf{z} \in D} \|f_{m+1}(\mathbf{z}), \dots, f_d(\mathbf{z})\|$ . Figure 4(b) shows the plot of the mean angular error against  $r^2$  for  $S^2 \times C$  in the noisy case. Figures 4(c)–(f) show the plots for  $S^5 \times C$  and  $S^8 \times C$ .

### 3.3 Curve

We experimented with the curve  $\varphi : [0, \pi] \rightarrow \mathbb{R}^d$  such that  $\varphi(\theta) = (\theta, \sin \theta, \dots, \sin \theta)$  for  $d = 20, 60$  and  $100$ . We pick the point  $\mathbf{p} = (\pi/2, 1, \dots, 1)$ , where the curve twists a lot. We vary

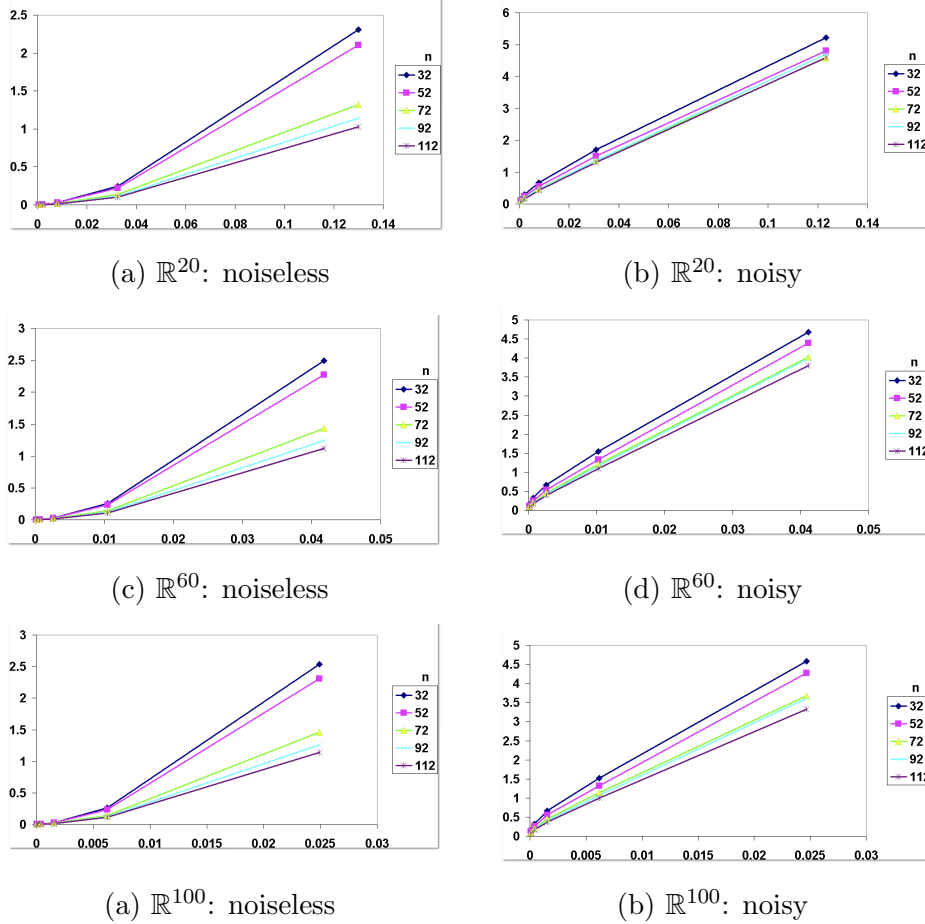


Figure 5: Plots of mean angular error for the sinusoidal curve in  $\mathbb{R}^{20}$ ,  $\mathbb{R}^{60}$  and  $\mathbb{R}^{100}$ . The vertical axes measure the mean angular error. The horizontal axes measure  $r^2/d$ . The maximum noise level is roughly 16%.

$r$  between 0 and  $\pi/2$ . For each  $r$ , the sample points are generated by sampling  $n$  values of  $\theta$  from  $[\pi/2 - r/\sqrt{d}, \pi/2 + r/\sqrt{d}]$ . The minimum value of  $\sin \theta$  over  $[\pi/2 - r/\sqrt{d}, \pi/2 + r/\sqrt{d}]$  is equal to  $\cos(r/\sqrt{d}) \approx 1 - r^2/(2d)$ . Therefore, the neighborhood radius is roughly  $(r^2/d + (d-1)r^4/(4d^2))^{1/2} \approx r/\sqrt{d}$ . We plot the mean angular error against  $r^2/20$  for  $\mathbb{R}^{20}$  in Figure 5(a).

Random noise is added by perturbing each sample point by a random displacement chosen from  $[0, 0.08r^2/\sqrt{d}]$  in a random direction in  $\mathbb{R}^d$ . As explained before, the minimum value of  $\sin \theta$  over  $[\pi/2 - r/\sqrt{d}, \pi/2 + r/\sqrt{d}]$  is roughly  $1 - r^2/(2d)$ . For every  $\ell \in [2, d]$ , the coordinate function  $f_\ell(\theta)$  is  $1 - \sin \theta \approx r^2/(2d)$ . So the maximum value of  $\|f_2(\theta), \dots, f_d(\theta)\|$  is roughly  $\sqrt{(d-1)r^4/(4d^2)} \approx r^2/(2\sqrt{d})$ . It follows that the maximum noise magnitude of  $0.08r^2/\sqrt{d}$  is roughly 16% of the maximum value of  $\|f_2(\theta), \dots, f_d(\theta)\|$ . Figure 5(b) shows the plot for the noisy cases in  $\mathbb{R}^{20}$ . Figures 5(c)–(f) show the plots for  $\mathbb{R}^{60}$  and  $\mathbb{R}^{100}$ .

## 4 Transformation to an eigenvalue problem

In this section, we show how to reduce our tangent estimation problem to finding the eigenvalues of the matrix  $\mathbf{H}\mathbf{H}^t$  as described in Section 2.2.

For every  $d$ -dimensional vector  $\hat{\mathbf{g}}_\ell$  that satisfies (2.5) and every  $\binom{d+1}{2}$ -dimensional vector  $\hat{\mathbf{c}}_\ell$ ,

we can apply (2.4) to evaluate the corresponding  $\widehat{F}_\ell$  at  $\mathbf{a}_p$  for  $p \in [1, n]$ . Since  $F_\ell(\mathbf{a}_p) = 0$  and  $\widehat{F}_\ell$  is supposed to approximate  $F_\ell$ , we should choose  $\widehat{\mathbf{g}}_\ell$  and  $\widehat{\mathbf{c}}_\ell$  to minimize some fitting error. Lemma 4.1 below shows that solving the eigenvalue problem for the matrix  $\mathbf{H}\mathbf{H}^t$  gives the best, less “curvy” fit.

Since  $F_\ell(\mathbf{a}_p) = 0$ , the mean squared error  $\frac{1}{n} \sum_{\ell=m+1}^d \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2$  is a natural error measure. However, since  $n \leq \binom{d+1}{2}$ , the system is under-determined and  $\widehat{F}_\ell(\mathbf{a}_p)$  can be made zero for all  $p \in [1, n]$  for many choices of  $\widehat{\mathbf{g}}_\ell$  and  $\widehat{\mathbf{c}}_\ell$ . We change the objective function to  $\frac{1}{n} \sum_{\ell=m+1}^d \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$  for some positive parameter  $\alpha$ . Intuitively,  $\widehat{\mathbf{c}}_\ell$  is an “approximation” of the linearization  $\mathbf{c}_\ell$  of  $\mathbf{Q}_\ell$ , so the penalty  $\sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$  favors a less “curvy” fit. We convert this optimization problem to a matrix problem as follows. Define:

$$\begin{aligned} d_0 &\stackrel{\text{def}}{=} \binom{d+1}{2} \\ \mathbf{U}_\alpha &\stackrel{\text{def}}{=} \frac{1}{n} \widehat{\mathbf{B}}^t \widehat{\mathbf{B}} + \alpha \mathbf{I}_{d_0} \\ \mathbf{H}_\alpha &\stackrel{\text{def}}{=} \frac{1}{n} \mathbf{A}^t \widehat{\mathbf{B}} \mathbf{U}_\alpha^{-1} \\ \mathbf{W}_\alpha &\stackrel{\text{def}}{=} \frac{1}{n} \mathbf{A}^t \mathbf{A} - \mathbf{H}_\alpha \mathbf{U}_\alpha \mathbf{H}_\alpha^t \end{aligned} \tag{4.1}$$

$\mathbf{U}_\alpha$  is a  $d_0 \times d_0$  matrix,  $\mathbf{H}_\alpha$  is an  $d \times d_0$  matrix, and  $\mathbf{W}_\alpha$  is a  $d \times d$  matrix.  $\mathbf{U}_\alpha$  and  $\mathbf{W}_\alpha$  are square and symmetric, and  $\mathbf{U}_\alpha$  is invertible because  $\alpha > 0$ . Lemma 4.1(i) below shows that the minimization of  $\frac{1}{n} \sum_{\ell=m+1}^d \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$  subject to (2.5) is equivalent to finding the smallest  $d - m$  eigenvalues of  $\mathbf{W}_\alpha$ . Lemma 4.1(i) also gives the optimal setting of  $\widehat{\mathbf{c}}_\ell$ .

A typical experimental approach is to solve the above optimization problem on some training data in order to tune the parameter  $\alpha$ . It may thus be time-consuming to set  $\alpha$  appropriately. For a small enough  $\alpha$ , the quantity  $\frac{1}{n} \sum_{\ell=m+1}^d \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2$  will be made zero as the system is under-determined, which means that the objective function value is effectively  $\sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$ . Thus, a curvy fit is penalized no matter how small  $\alpha$  is. This motivates us to push  $\alpha$  to zero in the limit. Lemma 4.1(ii) and (iii) show that  $\lim_{\alpha \rightarrow 0} \frac{1}{\alpha} \mathbf{W}_\alpha$  gives another eigenvalue problem that does not require any parameter, which is the problem solved by our tangent estimation algorithm in Section 2.2.

#### Lemma 4.1

- (i) For all  $\alpha > 0$  and for all mutually orthogonal unit vectors  $\widehat{\mathbf{g}}_{m+1}, \dots, \widehat{\mathbf{g}}_d$ , the value of  $\frac{1}{n} \sum_{\ell=m+1}^d \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$  is minimized with respect to  $\alpha$  and  $\widehat{\mathbf{g}}_{m+1}, \dots, \widehat{\mathbf{g}}_d$  when  $\widehat{\mathbf{c}}_\ell = -\mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell$  for every  $\ell \in [m+1, d]$ , and this minimum is equal to  $\sum_{\ell=m+1}^d \widehat{\mathbf{g}}_\ell^t \mathbf{W}_\alpha \widehat{\mathbf{g}}_\ell$ .
- (ii)  $\frac{1}{\alpha} \mathbf{W}_\alpha = \mathbf{A}^t \mathbf{L} \Sigma_\alpha \mathbf{L}^t \mathbf{A}$ , where  $\Sigma_\alpha = \text{diag}_n \left( \frac{1}{\alpha n + \lambda_1^2}, \dots, \frac{1}{\alpha n + \lambda_{m_0}^2}, \frac{1}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{1}{\alpha n + \lambda_{m_0+1}^2} \right)$ .
- (iii) Let  $\mathbf{H} = \lim_{\alpha \rightarrow 0} \mathbf{H}_\alpha$ . Let  $\Sigma = \lim_{\alpha \rightarrow 0} \Sigma_\alpha$ . Then,  $\mathbf{H} = (\widehat{\mathbf{B}}^t \mathbf{A})^t$  and  $\mathbf{A}^t \mathbf{L} \Sigma \mathbf{L}^t \mathbf{A} = \mathbf{H} \mathbf{H}^t$ . For every  $\ell \in [m+1, d]$ , if we set  $\widehat{\mathbf{g}}_\ell$  to be the unit eigenvector corresponding to the  $\ell$ th largest eigenvalue of  $\mathbf{A}^t \mathbf{L} \Sigma \mathbf{L}^t \mathbf{A}$  and set  $\widehat{\mathbf{c}}_\ell = -\mathbf{H}^t \widehat{\mathbf{g}}_\ell$ , then  $\widehat{F}_\ell(\mathbf{a}_p) = 0$  for every  $p \in [1, n]$ .

*Proof.* We follow the argument of Taubin who proved a result similar to (i) for reconstructing algebraic curves [32] in the presence of enough sample points so that the system is not under-



determined and the penalty  $\sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$  is not needed. By (2.4),

$$\begin{aligned} \frac{1}{n} \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 &= \frac{1}{n} (\widehat{\mathbf{g}}_\ell^t \ \widehat{\mathbf{c}}_\ell^t) \begin{pmatrix} \mathbf{A}^t \\ \widehat{\mathbf{B}}^t \end{pmatrix} (\mathbf{A} \ \widehat{\mathbf{B}}) \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell \end{pmatrix} \\ &= (\widehat{\mathbf{g}}_\ell^t \ \widehat{\mathbf{c}}_\ell^t) \begin{pmatrix} \frac{1}{n} \mathbf{A}^t \mathbf{A} & \frac{1}{n} \mathbf{A}^t \widehat{\mathbf{B}} \\ \frac{1}{n} \widehat{\mathbf{B}}^t \mathbf{A} & \frac{1}{n} \widehat{\mathbf{B}}^t \widehat{\mathbf{B}} \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell \end{pmatrix}. \end{aligned}$$

Also,  $\alpha \|\widehat{\mathbf{c}}_\ell\|^2 = (\widehat{\mathbf{g}}_\ell^t \ \widehat{\mathbf{c}}_\ell^t) \begin{pmatrix} \mathbf{0}_{d,d} & \mathbf{0}_{d,d_0} \\ \mathbf{0}_{d_0,d} & \alpha \mathbf{I}_{d_0} \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell \end{pmatrix}$ . Adding it to the above gives:

$$\frac{1}{n} \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \alpha \|\widehat{\mathbf{c}}_\ell\|^2 = (\widehat{\mathbf{g}}_\ell^t \ \widehat{\mathbf{c}}_\ell^t) \mathbf{Z} \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell \end{pmatrix},$$

where

$$\mathbf{Z} = \begin{pmatrix} \frac{1}{n} \mathbf{A}^t \mathbf{A} & \frac{1}{n} \mathbf{A}^t \widehat{\mathbf{B}} \\ \frac{1}{n} \widehat{\mathbf{B}}^t \mathbf{A} & \frac{1}{n} \widehat{\mathbf{B}}^t \widehat{\mathbf{B}} + \alpha \mathbf{I}_{d_0} \end{pmatrix} = \begin{pmatrix} \mathbf{W}_\alpha + \mathbf{H}_\alpha \mathbf{U}_\alpha \mathbf{H}_\alpha^t & \mathbf{H}_\alpha \mathbf{U}_\alpha \\ \mathbf{U}_\alpha \mathbf{H}_\alpha^t & \mathbf{U}_\alpha \end{pmatrix}.$$

We use the fact that  $\mathbf{U}_\alpha$  is symmetric in deriving  $\mathbf{U}_\alpha \mathbf{H}_\alpha^t$  in the lower left quadrant of  $\mathbf{Z}$ . We write  $\mathbf{Z}$  as the product of three matrices:

$$\mathbf{Z} = \begin{pmatrix} \mathbf{I}_d & \mathbf{H}_\alpha \\ \mathbf{0}_{d_0,d} & \mathbf{I}_{d_0} \end{pmatrix} \begin{pmatrix} \mathbf{W}_\alpha & \mathbf{0}_{d,d_0} \\ \mathbf{0}_{d_0,d} & \mathbf{U}_\alpha \end{pmatrix} \begin{pmatrix} \mathbf{I}_d & \mathbf{0}_{d,d_0} \\ \mathbf{H}_\alpha^t & \mathbf{I}_{d_0} \end{pmatrix}.$$

Then,

$$\begin{aligned} (\widehat{\mathbf{g}}_\ell^t \ \widehat{\mathbf{c}}_\ell^t) \mathbf{Z} \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell \end{pmatrix} &= (\widehat{\mathbf{g}}_\ell^t \ \widehat{\mathbf{c}}_\ell^t) \begin{pmatrix} \mathbf{I}_d & \mathbf{H}_\alpha \\ \mathbf{0}_{d_0,d} & \mathbf{I}_{d_0} \end{pmatrix} \begin{pmatrix} \mathbf{W}_\alpha & \mathbf{0}_{d,d_0} \\ \mathbf{0}_{d_0,d} & \mathbf{U}_\alpha \end{pmatrix} \begin{pmatrix} \mathbf{I}_d & \mathbf{0}_{d,d_0} \\ \mathbf{H}_\alpha^t & \mathbf{I}_{d_0} \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell \end{pmatrix} \\ &= (\widehat{\mathbf{g}}_\ell^t \ \widehat{\mathbf{c}}_\ell^t + \widehat{\mathbf{g}}_\ell^t \mathbf{H}_\alpha) \begin{pmatrix} \mathbf{W}_\alpha & \mathbf{0}_{d,d_0} \\ \mathbf{0}_{d_0,d} & \mathbf{U}_\alpha \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{g}}_\ell \\ \widehat{\mathbf{c}}_\ell + \mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell \end{pmatrix} \\ &= \widehat{\mathbf{g}}_\ell^t \mathbf{W}_\alpha \widehat{\mathbf{g}}_\ell + (\widehat{\mathbf{c}}_\ell + \mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell)^t \mathbf{U}_\alpha (\widehat{\mathbf{c}}_\ell + \mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell). \end{aligned}$$

Observe that  $\mathbf{U}_\alpha$  is positive semidefinite. It implies that  $(\widehat{\mathbf{c}}_\ell + \mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell)^t \mathbf{U}_\alpha (\widehat{\mathbf{c}}_\ell + \mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell) \geq 0$ . The minimum value of  $\frac{1}{n} \sum_{\ell=m+1}^d \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \sum_{\ell=m+1}^d \alpha \|\widehat{\mathbf{c}}_\ell\|^2$  is thus achieved when  $\widehat{\mathbf{c}}_\ell = -\mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell$  for every  $\ell \in [m+1, d]$ , and this minimum equals  $\sum_{\ell=m+1}^d \widehat{\mathbf{g}}_\ell^t \mathbf{W}_\alpha \widehat{\mathbf{g}}_\ell$ .

Consider (ii). Plugging the definition  $\mathbf{H}_\alpha = \frac{1}{n} \mathbf{A}^t \widehat{\mathbf{B}} \mathbf{U}_\alpha^{-1}$  into the definition of  $\mathbf{W}_\alpha$  gives:

$$\begin{aligned} \mathbf{W}_\alpha &= \frac{1}{n} \mathbf{A}^t \mathbf{A} - \mathbf{H}_\alpha \mathbf{U}_\alpha \mathbf{H}_\alpha^t \\ &= \frac{1}{n} \mathbf{A}^t \mathbf{A} - \frac{1}{n} \mathbf{A}^t \widehat{\mathbf{B}} \mathbf{H}_\alpha^t \\ &= \frac{1}{n} \mathbf{A}^t \mathbf{A} - \frac{1}{n^2} \mathbf{A}^t \widehat{\mathbf{B}} (\mathbf{U}_\alpha^{-1})^t \widehat{\mathbf{B}}^t \mathbf{A} \\ &= \frac{1}{n} \mathbf{A}^t \left( \mathbf{I}_n - \frac{1}{n} \widehat{\mathbf{B}} (\mathbf{U}_\alpha^{-1})^t \widehat{\mathbf{B}}^t \right) \mathbf{A}. \end{aligned} \tag{4.2}$$

Recall that  $m_0 < n \leq d_0$  by assumption and  $\widehat{\mathbf{L}} \widehat{\mathbf{A}} \widehat{\mathbf{R}}^t$  is the thin SVD of  $\widehat{\mathbf{B}}$ . Let  $\mathbf{L} (\widehat{\mathbf{L}} \ \mathbf{0}_{n,d_0-n}) \widehat{\mathbf{R}}^t$  be the full SVD of  $\widehat{\mathbf{B}}$ . So  $\mathbf{L}$  is an  $n \times n$  orthogonal matrix and  $\widehat{\mathbf{R}}$  is a  $d_0 \times d_0$  orthogonal matrix,

which means  $\mathbf{L}^t = \mathbf{L}^{-1}$  and  $\bar{\mathbf{R}}^t = \bar{\mathbf{R}}^{-1}$ . Moreover,  $\mathbf{R}$  is the leftmost  $d_0 \times n$  submatrix of  $\bar{\mathbf{R}}$ .

$$\begin{aligned}
\mathbf{U}_\alpha &= \frac{1}{n} \widehat{\mathbf{B}}^t \widehat{\mathbf{B}} + \alpha \mathbf{I}_{d_0} \\
&= \frac{1}{n} \bar{\mathbf{R}} \left( \widehat{\Lambda} \ \mathbf{0}_{n, d_0-n} \right)^t \left( \widehat{\Lambda} \ \mathbf{0}_{n, d_0-n} \right) \bar{\mathbf{R}}^t + \alpha \mathbf{I}_{d_0} \\
&= \bar{\mathbf{R}} \cdot \text{diag}_{d_0} \left( \frac{\lambda_1^2}{n}, \dots, \frac{\lambda_{m_0}^2}{n}, \frac{\lambda_{m_0+1}^2}{n}, \dots, \frac{\lambda_{m_0+1}^2}{n}, \underbrace{0, \dots, 0}_{d_0-n \text{ copies}} \right) \cdot \bar{\mathbf{R}}^t + \bar{\mathbf{R}}(\alpha \mathbf{I}_{d_0}) \bar{\mathbf{R}}^t \\
&= \bar{\mathbf{R}} \cdot \text{diag}_{d_0} \left( \frac{\alpha n + \lambda_1^2}{n}, \dots, \frac{\alpha n + \lambda_{m_0}^2}{n}, \frac{\alpha n + \lambda_{m_0+1}^2}{n}, \dots, \frac{\alpha n + \lambda_{m_0+1}^2}{n}, \underbrace{\alpha, \dots, \alpha}_{d_0-n \text{ copies}} \right) \cdot \bar{\mathbf{R}}^t.
\end{aligned}$$

It implies that

$$\mathbf{U}_\alpha^{-1} = \bar{\mathbf{R}} \cdot \text{diag}_{d_0} \left( \frac{n}{\alpha n + \lambda_1^2}, \dots, \frac{n}{\alpha n + \lambda_{m_0}^2}, \frac{n}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{n}{\alpha n + \lambda_{m_0+1}^2}, \frac{1}{\alpha}, \dots, \frac{1}{\alpha} \right) \cdot \bar{\mathbf{R}}^t. \quad (4.3)$$

Therefore,

$$\begin{aligned}
\frac{1}{n} \widehat{\mathbf{B}} (\mathbf{U}_\alpha^{-1})^t \widehat{\mathbf{B}}^t &= \mathbf{L} \left( \widehat{\Lambda} \ \mathbf{0}_{n, d_0-n} \right) \cdot \\
&\quad \text{diag}_{d_0} \left( \frac{1}{\alpha n + \lambda_1^2}, \dots, \frac{1}{\alpha n + \lambda_{m_0}^2}, \frac{1}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{1}{\alpha n + \lambda_{m_0+1}^2}, \frac{1}{\alpha n}, \dots, \frac{1}{\alpha n} \right) \cdot \\
&\quad \left( \widehat{\Lambda} \ \mathbf{0}_{n, d_0-n} \right)^t \mathbf{L}^t \\
&= \mathbf{L} \cdot \text{diag}_n \left( \frac{\lambda_1^2}{\alpha n + \lambda_1^2}, \dots, \frac{\lambda_{m_0}^2}{\alpha n + \lambda_{m_0}^2}, \frac{\lambda_{m_0+1}^2}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{\lambda_{m_0+1}^2}{\alpha n + \lambda_{m_0+1}^2} \right) \cdot \mathbf{L}^t.
\end{aligned}$$

Plugging this into (4.2) gives:

$$\begin{aligned}
\mathbf{W}_\alpha &= \frac{1}{n} \mathbf{A}^t \mathbf{L} \left( \mathbf{I}_n - \text{diag}_n \left( \frac{\lambda_1^2}{\alpha n + \lambda_1^2}, \dots, \frac{\lambda_{m_0}^2}{\alpha n + \lambda_{m_0}^2}, \frac{\lambda_{m_0+1}^2}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{\lambda_{m_0+1}^2}{\alpha n + \lambda_{m_0+1}^2} \right) \right) \mathbf{L}^t \mathbf{A} \\
&= \frac{1}{n} \mathbf{A}^t \mathbf{L} \cdot \text{diag}_n \left( \frac{\alpha n}{\alpha n + \lambda_1^2}, \dots, \frac{\alpha n}{\alpha n + \lambda_{m_0}^2}, \frac{\alpha n}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{\alpha n}{\alpha n + \lambda_{m_0+1}^2} \right) \cdot \mathbf{L}^t \mathbf{A} \\
&= \alpha \mathbf{A}^t \mathbf{L} \Sigma_\alpha \mathbf{L}^t \mathbf{A}.
\end{aligned}$$

Consider (iii). Using the definition of  $\mathbf{H}_\alpha$ , (4.3), and the full SVD of  $\widehat{\mathbf{B}}$ , we obtain

$$\begin{aligned}
\mathbf{H}_\alpha &= \frac{1}{n} \mathbf{A}^t \widehat{\mathbf{B}} \mathbf{U}_\alpha^{-1} \\
&= \mathbf{A}^t \mathbf{L} \left( \widehat{\Lambda} \ \mathbf{0}_{n, d_0-n} \right) \cdot \text{diag}_{d_0} \left( \frac{1}{\alpha n + \lambda_1^2}, \dots, \frac{1}{\alpha n + \lambda_{m_0}^2}, \frac{1}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{1}{\alpha n + \lambda_{m_0+1}^2}, \frac{1}{\alpha n}, \dots, \frac{1}{\alpha n} \right) \cdot \bar{\mathbf{R}}^t \\
&= \mathbf{A}^t \mathbf{L} \cdot \text{diag}_n \left( \frac{\lambda_1}{\alpha n + \lambda_1^2}, \dots, \frac{\lambda_{m_0}}{\alpha n + \lambda_{m_0}^2}, \frac{\lambda_{m_0+1}}{\alpha n + \lambda_{m_0+1}^2}, \dots, \frac{\lambda_{m_0+1}}{\alpha n + \lambda_{m_0+1}^2} \right) \cdot \mathbf{R}^t.
\end{aligned}$$

In the last step, we use the property that  $\mathbf{R}$  is the leftmost  $d_0 \times n$  submatrix of  $\bar{\mathbf{R}}$ . Define  $\mathbf{H} = \lim_{\alpha \rightarrow 0} \mathbf{H}_\alpha$  and  $\Sigma = \lim_{\alpha \rightarrow 0} \Sigma_\alpha$ . As a result,

$$\mathbf{H} = \lim_{\alpha \rightarrow 0} \mathbf{H}_\alpha = \mathbf{A}^t \mathbf{L} \cdot \text{diag}_n \left( \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_{m_0}}, \frac{1}{\lambda_{m_0+1}}, \dots, \frac{1}{\lambda_{m_0+1}} \right) \cdot \mathbf{R}^t.$$

Observe that  $\widehat{\mathbf{B}}^\dagger = \mathbf{R} \cdot \text{diag}_n \left( \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_{m_0}}, \frac{1}{\lambda_{m_0+1}}, \dots, \frac{1}{\lambda_{m_0+1}} \right) \cdot \mathbf{L}^t$ . Therefore,  $\mathbf{H} = (\widehat{\mathbf{B}}^\dagger \mathbf{A})^t$  and  $\mathbf{H} \mathbf{H}^t = \mathbf{A}^t \mathbf{L} \cdot \text{diag}_n \left( \frac{1}{\lambda_1^2}, \dots, \frac{1}{\lambda_{m_0}^2}, \frac{1}{\lambda_{m_0+1}^2}, \dots, \frac{1}{\lambda_{m_0+1}^2} \right) \cdot \mathbf{L}^t \mathbf{A} = \mathbf{A}^t \mathbf{L} \Sigma \mathbf{L}^t \mathbf{A}$ .

Take any  $\ell \in [m+1, d]$ . Let  $\mathbf{e}_\alpha$  be the unit eigenvector corresponding to the  $\ell$ th largest eigenvalue of  $\mathbf{A}^t \mathbf{L} \Sigma_\alpha \mathbf{L}^t \mathbf{A}$  for a positive but arbitrarily small  $\alpha$ . By (i),  $\widehat{\mathbf{c}}_\ell$  should be set to

$-\mathbf{H}_\alpha^t \mathbf{e}_\alpha$  to minimize the fitting error, which makes  $\mathbf{e}_\alpha^t \mathbf{W}_\alpha \mathbf{e}_\alpha = \frac{1}{n} \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \alpha \|\widehat{\mathbf{c}}_\ell\|^2 = \frac{1}{n} \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \alpha \mathbf{e}_\alpha^t (\mathbf{H}_\alpha \mathbf{H}_\alpha^t) \mathbf{e}_\alpha$ . By (ii),

$$\mathbf{e}_\alpha^t (\mathbf{A}^t \mathbf{L} \Sigma_\alpha \mathbf{L}^t \mathbf{A}) \mathbf{e}_\alpha = \mathbf{e}_\alpha^t (\frac{1}{\alpha} \mathbf{W}_\alpha) \mathbf{e}_\alpha = \frac{1}{\alpha n} \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 + \mathbf{e}_\alpha^t (\mathbf{H}_\alpha \mathbf{H}_\alpha^t) \mathbf{e}_\alpha.$$

As  $\alpha \rightarrow 0$ , the left hand side and the second term on the right hand side converge to the same value. Therefore, if we set  $\widehat{\mathbf{g}}_\ell = \mathbf{e}_\alpha$  and  $\widehat{\mathbf{c}}_\ell = -\mathbf{H}_\alpha^t \widehat{\mathbf{g}}_\ell$ , then  $\lim_{\alpha \rightarrow 0} \frac{1}{\alpha n} \sum_{p=1}^n \widehat{F}_\ell(\mathbf{a}_p)^2 = 0$ , which makes  $\lim_{\alpha \rightarrow 0} \widehat{F}_\ell(\mathbf{a}_p) = 0$  for every  $p \in [1, n]$ . This proves (iii).  $\square$

**Remark:** The condition  $n \leq \binom{d+1}{2}$  is needed for the proofs of Lemma 4.1(ii) and the identity  $\mathbf{A}^t \mathbf{L} \Sigma \mathbf{L}^t \mathbf{A} = \mathbf{H} \mathbf{H}^t$  in Lemma 4.1(iii) to go through. The subsequent error analysis studies the eigenvalues of  $\mathbf{H} \mathbf{H}^t$  and use the identity  $\mathbf{A}^t \mathbf{L} \Sigma \mathbf{L}^t \mathbf{A} = \mathbf{H} \mathbf{H}^t$  to bound the angular error.

## 5 Derivatives and Sums

### 5.1 Derivatives of coordinate functions

We derive two results in this subsection. Lemma 5.1(i) puts an upper bound on the 2-norm of the second derivatives of the coordinate functions. Lemma 5.1(ii) bounds the error of approximating the coordinate functions using their second derivatives. We use Taylor expansions [23] heavily in the proofs. Given a smooth function  $h : \mathcal{X} \rightarrow \mathbb{R}$ , where  $\mathcal{X}$  is an open subset of  $\mathbb{R}^k$  for some  $k \geq 1$ , if two points  $\mathbf{a}, \mathbf{x} \in \mathcal{X}$  are connected by a segment in  $\mathcal{X}$ , then for every integer  $j \geq 1$ , there exists a point  $\mathbf{z}$  in the interior of the segment connecting  $\mathbf{a}$  and  $\mathbf{x}$  such that

$$\begin{aligned} h(\mathbf{x}) &= h(\mathbf{a}) + \mathbf{D}h|_{\mathbf{a}}(\mathbf{x} - \mathbf{a}) + \cdots + \frac{1}{(j-1)!} \mathbf{D}^{j-1}h|_{\mathbf{a}}(\mathbf{x} - \mathbf{a}, \dots, \mathbf{x} - \mathbf{a}) + \\ &\quad \frac{1}{j!} \mathbf{D}^j h|_{\mathbf{z}}(\mathbf{x} - \mathbf{a}, \dots, \mathbf{x} - \mathbf{a}). \end{aligned}$$

If we let  $\mathbf{u}$  be the unit vector  $(\mathbf{x} - \mathbf{a})/\|\mathbf{x} - \mathbf{a}\|$  and let  $\varepsilon = \|\mathbf{x} - \mathbf{a}\|$ , then

$$h(\mathbf{x}) = h(\mathbf{a}) + \varepsilon \mathbf{D}h|_{\mathbf{a}}(\mathbf{u}) + \cdots + \frac{\varepsilon^{j-1}}{(j-1)!} \mathbf{D}^{j-1}h|_{\mathbf{a}}(\mathbf{u}, \dots, \mathbf{u}) + \frac{\varepsilon^j}{j!} \mathbf{D}^j h|_{\mathbf{z}}(\mathbf{u}, \dots, \mathbf{u}).$$

As  $\varepsilon = \|\mathbf{x} - \mathbf{a}\|$  approaches zero, the magnitude of the remainder  $\frac{\varepsilon^j}{j!} \mathbf{D}^j h|_{\mathbf{z}}(\mathbf{u}, \dots, \mathbf{u})$  approaches zero faster than the magnitude of every preceding term in the expansion. In our case, we will use the Taylor expansion of a coordinate function  $f_\ell$  around the origin, and there are no constant and linear term in the expansion.

**Lemma 5.1** *There exists a value  $\varepsilon_0 \in (0, 1)$  such that for every  $\varepsilon \in (0, \varepsilon_0]$  and every unit vector  $\mathbf{u} \in \mathbb{R}^m$ , the following properties hold.*

(i)  $\|(\mathbf{D}^2 f_{m+1}|_0(\mathbf{u}, \mathbf{u}) \cdots \mathbf{D}^2 f_d|_0(\mathbf{u}, \mathbf{u}))\| = O(1/\varrho)$ .

(ii) Let  $\mathbf{v} = \varrho \varepsilon \mathbf{u}$ .  $\|(f_{m+1}(\mathbf{v}) - \frac{1}{2} \mathbf{D}^2 f_{m+1}|_0(\mathbf{v}, \mathbf{v}) \cdots f_d(\mathbf{v}) - \frac{1}{2} \mathbf{D}^2 f_d|_0(\mathbf{v}, \mathbf{v}))\| = O(\varrho \varepsilon^3)$ .

*Proof.* Rotate space such that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$  and the coordinate axes  $x_{m+1}, \dots, x_d$  span the normal space of  $\mathcal{M}$  at the origin. Let  $\mathbf{v} = \varrho \varepsilon \mathbf{u}$ .

Consider (i). Take any  $\ell \in [m+1, d]$ . We claim that if  $\varepsilon_0$  is sufficiently small, then  $|f_\ell(\mathbf{v})| \geq \frac{1}{4} \varrho^2 \varepsilon^2 |\mathbf{D}^2 f_\ell|_0(\mathbf{u}, \mathbf{u})|$  for all  $\varepsilon \in (0, \varepsilon_0]$ . If  $\mathbf{D}^2 f_\ell|_0(\mathbf{u}, \mathbf{u}) = 0$ , our claim is trivially true. Assume that

$D^2f_\ell|_0(\mathbf{u}, \mathbf{u}) \neq 0$ . The Taylor expansion of  $f_\ell(\mathbf{v})$  is  $\frac{1}{2}D^2f_\ell|_0(\mathbf{v}, \mathbf{v}) + \dots = \frac{1}{2}\varrho^2\varepsilon^2D^2f_\ell|_0(\mathbf{u}, \mathbf{u}) + \dots$ . As  $\varepsilon_0$  approaches 0, the remainder approaches zero faster than  $\frac{1}{2}D^2f_\ell|_0(\mathbf{v}, \mathbf{v})$  does. Therefore, there exists  $\varepsilon_0 > 0$  such that for every  $\varepsilon \in (0, \varepsilon_0]$  and  $\mathbf{v} = \varrho\varepsilon\mathbf{u}$ ,

$$\left| f_\ell(\mathbf{v}) - \frac{1}{2}D^2f_\ell|_0(\mathbf{v}, \mathbf{v}) \right| \leq \frac{1}{4} |D^2f_\ell|_0(\mathbf{v}, \mathbf{v})| = \frac{1}{4}\varrho^2\varepsilon^2 |D^2f_\ell|_0(\mathbf{u}, \mathbf{u})|.$$

Thus,  $|f_\ell(\mathbf{v})| \geq \frac{1}{4}\varrho^2\varepsilon^2 |D^2f_\ell|_0(\mathbf{u}, \mathbf{u})|$ , proving our claim.

Any point in  $\mathcal{M}$  within a distance  $\varrho\varepsilon$  from the origin is at distance  $O(\varrho\varepsilon^2)$  from  $\mathcal{T}$  [15, Lemma 6], which implies that  $\|(f_{m+1}(\mathbf{v}) \ \dots \ f_d(\mathbf{v}))\| = O(\varrho\varepsilon^2)$ . Therefore, by our claim,

$$\frac{1}{4}\varrho^2\varepsilon^2 \|(D^2f_{m+1}|_0(\mathbf{u}, \mathbf{u}) \ \dots \ D^2f_d|_0(\mathbf{u}, \mathbf{u}))\| \leq \|(f_{m+1}(\mathbf{v}) \ \dots \ f_d(\mathbf{v}))\| = O(\varrho\varepsilon^2).$$

Dividing both sides by  $\varrho^2\varepsilon^2/4$  gives  $\|(D^2f_{m+1}|_0(\mathbf{u}, \mathbf{u}) \ \dots \ D^2f_d|_0(\mathbf{u}, \mathbf{u}))\| = O(1/\varrho)$ , establishing the correctness of (i).

Consider (ii). Let  $B_\varepsilon$  denote the  $d$ -ball centered at the origin with radius  $\varrho\varepsilon$ . Take any  $\ell \in [m+1, d]$ . Define the plane  $\mathcal{L} = \{a\mathbf{u} + \mathbf{z} : a \in \mathbb{R} \wedge \mathbf{z} \text{ is a vector parallel to the } x_\ell \text{ axis}\}$ . Since  $\mathbf{u}$  is a vector in the tangent space of  $\mathcal{M}$  at the origin, for a sufficiently small  $\varepsilon_0$ ,  $\mathcal{M} \cap B_\varepsilon \cap \mathcal{L}$  is a one-dimensional curve in the plane  $\mathcal{L}$  and  $D^2f_\ell|_0(\mathbf{u}, \mathbf{u})$  is the rate of change of the  $\ell$ -th coordinate of the unit tangent along  $\mathcal{M} \cap B_\varepsilon \cap \mathcal{L}$  at the origin [28, Chapter 5]. That is,  $D^2f_\ell|_0(\mathbf{u}, \mathbf{u})$  is the reciprocal of the radius of curvature of  $\mathcal{M} \cap B_\varepsilon \cap \mathcal{L}$  at the origin.

If we scale down the unit length linearly, all lengths increase linearly; in particular, the radius of curvature of  $\mathcal{M} \cap B_\varepsilon \cap \mathcal{L}$  at the origin, the local feature size  $\varrho$  at the origin, and the value  $|f_\ell(\mathbf{v})|$ . The linear increase in the radius of curvature means that  $|D^2f_\ell|_0(\mathbf{u}, \mathbf{u})|$  decreases linearly, which implies that  $D^2f_\ell|_0(\mathbf{v}, \mathbf{v}) = \varrho^2\varepsilon^2D^2f_\ell|_0(\mathbf{u}, \mathbf{u})$  increases linearly.

A Taylor expansion of  $f_\ell(\mathbf{v})$  is  $\frac{1}{2}D^2f_\ell|_0(\mathbf{v}, \mathbf{v}) + \frac{1}{6}D^3f_\ell|_{\mathbf{p}_{\ell, \mathbf{u}, \varepsilon}}(\mathbf{v}, \mathbf{v}, \mathbf{v})$ , where  $\frac{1}{6}D^3f_\ell|_{\mathbf{p}_{\ell, \mathbf{u}, \varepsilon}}(\mathbf{v}, \mathbf{v}, \mathbf{v})$  is the remainder and  $\mathbf{p}_{\ell, \mathbf{u}, \varepsilon} = \varrho\omega_{\ell, \mathbf{u}, \varepsilon}\mathbf{u}$  for some  $\omega_{\ell, \mathbf{u}, \varepsilon} \in (0, \varepsilon)$ . Notice that  $\mathbf{p}_{\ell, \mathbf{u}, \varepsilon}$  and  $\omega_{\ell, \mathbf{u}, \varepsilon}$  depend on  $\varepsilon$  and  $\mathbf{u}$ . The remainder can be rewritten as:

$$f_\ell(\mathbf{v}) - \frac{1}{2}D^2f_\ell|_0(\mathbf{v}, \mathbf{v}) = \frac{1}{6}D^3f_\ell|_{\mathbf{p}_{\ell, \mathbf{u}, \varepsilon}}(\mathbf{v}, \mathbf{v}, \mathbf{v}) = \varrho\varepsilon^3 \cdot \frac{1}{6}D^3f_\ell|_{\mathbf{p}_{\ell, \mathbf{u}, \varepsilon}}(\mathbf{u}, \mathbf{u}, \mathbf{u}) \quad (5.1)$$

Either  $f_\ell(\mathbf{v}) - \frac{1}{2}D^2f_\ell|_0(\mathbf{v}, \mathbf{v})$  is zero or it changes linearly when we scale down the unit length because  $f_\ell(\mathbf{v})$  and  $\frac{1}{2}D^2f_\ell|_0(\mathbf{v}, \mathbf{v})$  do. The leading factor  $\varrho\varepsilon^3$  on the right hand side of (5.1) changes linearly too as we scale down the unit length. Thus,  $\frac{1}{6}\varrho^2D^3f_\ell|_{\mathbf{p}_{\ell, \mathbf{u}, \varepsilon}}(\mathbf{u}, \mathbf{u}, \mathbf{u})$  is scale independent.

We have the following equation:

$$\left\| \begin{pmatrix} f_{m+1}(\mathbf{v}) - \frac{1}{2}D^2f_{m+1}|_0(\mathbf{v}, \mathbf{v}) \\ \vdots \\ f_d(\mathbf{v}) - \frac{1}{2}D^2f_d|_0(\mathbf{v}, \mathbf{v}) \end{pmatrix} \right\| = \varrho\varepsilon^3 \cdot \left\| \begin{pmatrix} \frac{1}{6}\varrho^2D^3f_{m+1}|_{\mathbf{p}_{m+1, \mathbf{u}, \varepsilon}}(\mathbf{u}, \mathbf{u}, \mathbf{u}) \\ \vdots \\ \frac{1}{6}\varrho^2D^3f_d|_{\mathbf{p}_{d, \mathbf{u}, \varepsilon}}(\mathbf{u}, \mathbf{u}, \mathbf{u}) \end{pmatrix} \right\| \quad (5.2)$$

We maximize the norm of the vector on the right hand side of (5.2) by going over all possible unit vectors  $\mathbf{u}$  in the tangent space at the origin, but  $\mathbf{v}$  on the left hand side of (5.2) is kept fixed independent of the variation of  $\mathbf{u}$ . For every  $\mathbf{u}$ , we also vary  $\mathbf{p}_{\ell, \mathbf{u}, \varepsilon}$  by varying  $\varepsilon$  over  $(0, \varepsilon_0]$  to maximize the norm of the vector on the right hand side of (5.2). Therefore,

$$\left\| \begin{pmatrix} f_{m+1}(\mathbf{v}) - \frac{1}{2}D^2f_{m+1}|_0(\mathbf{v}, \mathbf{v}) \\ \vdots \\ f_d(\mathbf{v}) - \frac{1}{2}D^2f_d|_0(\mathbf{v}, \mathbf{v}) \end{pmatrix} \right\| \leq \varrho\varepsilon^3 \cdot \sup_{\substack{\|\mathbf{u}\|=1 \\ \varepsilon \in (0, \varepsilon_0]}} \left\| \begin{pmatrix} \frac{1}{6}\varrho^2D^3f_{m+1}|_{\mathbf{p}_{m+1, \mathbf{u}, \varepsilon}}(\mathbf{u}, \mathbf{u}, \mathbf{u}) \\ \vdots \\ \frac{1}{6}\varrho^2D^3f_d|_{\mathbf{p}_{d, \mathbf{u}, \varepsilon}}(\mathbf{u}, \mathbf{u}, \mathbf{u}) \end{pmatrix} \right\|.$$

In the inequality above, the rightmost factor on the right hand side is a constant that depends on  $\mathcal{M}$  but not on scale. This establishes the correctness of (ii).  $\square$

## 5.2 Sums of powers of coordinates

The main result in this subsection is Lemma 5.4 which bounds certain sums of powers of sample point coordinates. These bounds are proved using integration, the Chebyshev's inequality, and the technical result stated in Lemma 5.3 below. We need some notation. Let  $p$  and  $q$  denote two non-negative integers.

$$\beta(q) \stackrel{\text{def}}{=} \begin{cases} \frac{\pi(q-1)(q-3)\cdots 1}{q(q-2)\cdots 2}, & \text{if } q \text{ is even and positive,} \\ \frac{2(q-1)(q-3)\cdots 2}{q(q-2)\cdots 1}, & \text{if } q \text{ is odd.} \end{cases}$$

$$V_q \stackrel{\text{def}}{=} \text{volume of a unit } q\text{-ball.}$$

Note that  $\beta(1) = 2$  and  $V_0 = 1$ . It can be verified that  $\beta$  and  $V_q$  satisfy the following recurrences:

$$\forall q \geq 2, \quad \beta(q-1)\beta(q+2) = 2\pi(q+1)/(q(q+2)) \quad (5.3)$$

$$\forall q \geq 1, \quad V_q = \beta(q)V_{q-1} \quad (5.4)$$

$$\forall q \geq 2, \quad V_q = 2\pi V_{q-2}/q \quad (5.5)$$

We will need the following two technical results. The Chebyshev's inequality bounds the probability of a random variable deviating from the mean by a multiple of the standard deviation.

**Lemma 5.2 (Chebyshev's Inequality)** *Let  $Y$  be a random variable with finite expected value  $\mu$  and finite positive variance  $\sigma^2$ . For any positive real number  $a$ ,  $\Pr(|Y - \mu| \geq a\sigma) \leq 1/a^2$ .*

The next result gives the value of a particular integral that will be used often in the proof of Lemma 5.4.

**Lemma 5.3** *Let  $p$  and  $q$  be two non-negative integers. Let  $r_i$ ,  $x_i$  and  $r_{i+1}$  be three variables such that  $x_i \in [-r_i, r_i]$  and  $r_{i+1}^2 = r_i^2 - x_i^2$ . If  $p$  is non-negative and even, then*

$$\int_{-r_i}^{r_i} x_i^p r_{i+1}^q dx_i = \frac{(p-1)(p-3)\cdots 1}{(p+q+1)(p+q-1)\cdots (q+3)} \cdot \beta(q+1) \cdot r_i^{p+q+1},$$

where  $\frac{(p-1)(p-3)\cdots 1}{(p+q+1)(p+q-1)\cdots (q+3)}$  is interpreted as 1 when  $p = 0$ .

*Proof.* Perform a change of variables:  $x_i = r_i \sin \theta$  and  $r_{i+1} = r_i \cos \theta$  and we obtain

$$\int x_i^p r_{i+1}^q dx_i = r_i^{p+q+1} \int \sin^p \theta \cos^{q+1} \theta d\theta.$$

The limits of the integrals also change:  $[-r_i, r_i]$  becomes  $[-\pi/2, \pi/2]$ . The following two recursive formulae are from [19]:

$$\forall p \geq 2, q \geq 0, \quad \int \sin^p \theta \cos^{q+1} \theta d\theta = -\frac{\sin^{p-1} \theta \cos^{q+2} \theta}{p+q+1} + \frac{p-1}{p+q+1} \cdot \int \sin^{p-2} \theta \cos^{q+1} \theta d\theta \quad (5.6)$$

$$\forall p \geq 0, q \geq 1, \quad \int \sin^p \theta \cos^{q+1} \theta d\theta = \frac{\sin^{p+1} \theta \cos^q \theta}{p+q+1} + \frac{q}{p+q+1} \cdot \int \sin^p \theta \cos^{q-1} \theta d\theta. \quad (5.7)$$

In our case, since the limits of the integrals are  $[-\pi/2, \pi/2]$ , the leading additive terms in both (5.6) and (5.7) vanish. Then, we first repeatedly apply (5.6) to decrease the exponent  $p$ . Since  $p$  is even, it will eventually become zero and we then repeatedly apply (5.7) to decrease the exponent  $q$ . One can then verify that we obtain the result stated in the lemma.  $\square$

Lemma 5.4 below gives upper and lower bounds on certain sums of powers of sample point coordinates. These upper and lower bounds are obtained by calculating the expected values and variances of the sums of powers. Then, the Chebyshev's inequality is applied to obtain the high probability bound. Recall that the sample points in  $\mathcal{M}$  are generated by a Poisson process and the neighborhood of the origin being examined contains exactly  $n$  sample points (excluding the origin). It follows that these  $n$  sample points are uniformly distributed in that neighborhood [5]. As a result, the probability of a sample point falling into a region within the neighborhood is the ratio of the volume of that region to the neighborhood volume. This observation allows us to calculate the expected values and variances of the sums of powers by integration and applying Lemma 5.3.

**Lemma 5.4** *Assume that the coordinate axes  $x_1, \dots, x_m$  span the tangent space  $\mathcal{T}$  of  $\mathcal{M}$  at the origin. Let  $i, j, k$ , and  $l$  be four distinct integers from  $[1, m]$ . If  $\varepsilon$  is sufficiently small, then for every constant  $c > 0$ , the following properties hold simultaneously with probability  $1 - O(n^{-1/3})$ , where the hidden constant in the probability bound depends on  $c$ .*

- (i)  $\left| \sum_{p=1}^n a_{pi}^4 - \frac{3}{(m+2)(m+4)} n \varrho^4 \varepsilon^4 \right| < cn^{2/3} \varrho^4 \varepsilon^4 + O(n \varrho^4 \varepsilon^6)$ .
- (ii)  $\left| \sum_{p=1}^n a_{pi}^2 a_{pj}^2 - \frac{1}{(m+2)(m+4)} n \varrho^4 \varepsilon^4 \right| < cn^{2/3} \varrho^4 \varepsilon^4 + O(n \varrho^4 \varepsilon^6)$ .
- (iii)  $\left| \sum_{p=1}^n a_{pi}^2 - \frac{1}{m+2} n \varrho^2 \varepsilon^2 \right| < cn^{2/3} \varrho^2 \varepsilon^2 + O(n \varrho^2 \varepsilon^4)$ .
- (iv)  $\left| \sum_{p=1}^n a_{pi}^3 a_{pj} \right|, \left| \sum_{p=1}^n a_{pi}^2 a_{pj} a_{pk} \right|$  and  $\left| \sum_{p=1}^n a_{pi} a_{pj} a_{pk} a_{pl} \right|$  are less than  $cn^{2/3} \varrho^4 \varepsilon^4 + O(n \varrho^4 \varepsilon^6)$ .
- (v)  $\left| \sum_{p=1}^n a_{pi}^3 \right|, \left| \sum_{p=1}^n a_{pi}^2 a_{pj} \right|$  and  $\left| \sum_{p=1}^n a_{pi} a_{pj} a_{pk} \right|$  are less than  $cn^{2/3} \varrho^3 \varepsilon^3 + O(n \varrho^3 \varepsilon^5)$ .
- (vi)  $\left| \sum_{p=1}^n a_{pi} a_{pj} \right| < cn^{2/3} \varrho^2 \varepsilon^2 + O(n \varrho^2 \varepsilon^4)$ .

*Proof.* Let  $B_\varepsilon$  be the  $d$ -ball centered at the origin with radius  $\varrho\varepsilon$ . Since  $\varrho\varepsilon$  is defined to be the distance from the origin to the  $(n+1)$ -th sample point, there are exactly  $n$  sample points in the interior of  $B_\varepsilon$ , excluding the origin. As the sample points are generated by a Poisson process, these  $n$  sample points are uniformly distributed in the interior of  $\mathcal{M} \cap B_\varepsilon$ .

Consider an  $m$ -dimensional triangulated hyperrectangle  $R$  in  $\mathcal{T} \cap B_\varepsilon$  with infinitesimal side lengths  $dx_1, \dots, dx_m$ . Suppose that  $R$  lies well inside  $B_\varepsilon$  so that its  $2^m$  vertices are the orthogonal projections of  $2^m$  points in  $\mathcal{M} \cap B_\varepsilon$  onto  $\mathcal{T}$ . Connect these  $2^m$  points in  $\mathcal{M} \cap B_\varepsilon$  as in the triangulation of  $R$  to produce some  $m$ -simplices. The union of these  $m$ -simplices is an infinitesimal volume  $U$  which is a deformed version of  $R$ . Let  $dV$  denote the volume of  $U$ . Since  $U$  projects orthogonally onto  $R$ , the volume of  $U$  is at least the volume of  $R$ , and therefore,  $dV \geq dx_1 \cdots dx_m$ . For every segment  $s'$  in an  $m$ -simplex in  $U$ , it projects to a segment  $s$  in an  $m$ -simplex in the triangulation of  $R$  such that  $\text{length}(s') = \sec \theta \cdot \text{length}(s)$ , where  $\theta$  is the acute angle between the support lines of  $s'$  and  $s$ . The angle  $\theta$  is no more than the angle between  $\mathcal{T}$  and the tangent space at an endpoint of  $s'$ , which is known to be  $O(\varepsilon)$  [9, Lemma 15]. Therefore,  $dV \leq (1 + O(\varepsilon^2))^m dx_1 \cdots dx_m = (1 + O(\varepsilon^2)) dx_1 \cdots dx_m$ . In the last step, we use the fact that  $m$  is a constant, so  $m \cdot O(\varepsilon^2) + \binom{m}{2} \cdot O(\varepsilon^4) + \cdots = O(\varepsilon^2)$ . It follows that:

$$dx_1 \cdots dx_m \leq dV \leq (1 + O(\varepsilon^2)) \cdot dx_1 \cdots dx_m \quad (5.8)$$

Define the following symbols:

$$\begin{aligned} r_0 &\stackrel{\text{def}}{=} \varrho\varepsilon \\ r_1 &\stackrel{\text{def}}{=} \kappa r_0 \text{ for some parameter } \kappa \text{ to be specified later} \\ r_{i+1}^2 &\stackrel{\text{def}}{=} r_i^2 - x_i^2, \text{ for every } i \in [1, m-1] \end{aligned}$$

Since the tangent space at any point in  $\mathcal{M} \cap B_\varepsilon$  makes an  $O(\varepsilon)$  angle with  $\mathcal{T}$  [9], the projection of  $\mathcal{M} \cap B_\varepsilon$  onto  $\mathcal{T}$  covers an  $m$ -ball centered at the origin with radius  $r_0 \cos(O(\varepsilon)) = (1 - O(\varepsilon^2))r_0$ . The volume of this  $m$ -ball is thus  $(1 - O(\varepsilon^2))^m V_m r_0^m \geq (1 - O(\varepsilon^2))V_m r_0^m$ . Combining with (5.8), we obtain:

$$(1 - O(\varepsilon^2))V_m r_0^m \leq \text{vol}(\mathcal{M} \cap B_\varepsilon) \leq (1 + O(\varepsilon^2))V_m r_0^m. \quad (5.9)$$

We will prove that each of (i)–(vi) holds with probability  $1 - O(n^{-1/3})$ . Therefore, they hold simultaneously with probability at least  $1 - 6 \cdot O(n^{-1/3}) = 1 - O(n^{-1/3})$  as well. In the following,  $\mathbf{z} = (z_1, z_2, \dots, z_d)^t$  denotes a random point in  $\mathcal{M} \cap B_\varepsilon$ .

*Analysis of (i):* The variance of  $\sum_{p=1}^n a_{pi}^4$  equals  $n\mathbf{E}[z_i^8] - n(\mathbf{E}[z_i^4])^2 \leq n\mathbf{E}[z_i^8] \leq n\varrho^8\varepsilon^8$  because  $|z_i| \leq \varrho\varepsilon$ . Lemma 5.2 implies that  $\left| \sum_{p=1}^n a_{pi}^4 - \mathbf{E} \left[ \sum_{p=1}^n a_{pi}^4 \right] \right| < cn^{2/3}\varrho^4\varepsilon^4$  with probability  $1 - c^{-2}n^{-1/3}$ . It remains to bound  $\mathbf{E} \left[ \sum_{p=1}^n a_{pi}^4 \right]$ .

Since  $\mathbf{E}[\sum_{p=1}^n a_{pi}^4] = \frac{n}{\text{vol}(\mathcal{M} \cap B_\varepsilon)} \int_{\mathcal{M} \cap B_\varepsilon} z_i^4 dV$ , it follows from (5.9) that:

$$\frac{(1 - O(\varepsilon^2))n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^4 dV \leq \mathbf{E} \left[ \sum_{p=1}^n a_{pi}^4 \right] \leq \frac{(1 + O(\varepsilon^2))n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^4 dV \quad (5.10)$$

We first calculate  $\frac{n}{V_m r_0^m} \int_D x_i^4 dx_m \cdots dx_1$ , where  $D$  is the  $m$ -ball in  $\mathcal{T}$  centered at the origin with radius  $r_1$ . Consider an  $(m - i + 1)$ -ball with radius  $r_i$  centered at the origin. Its volume, which is  $V_{m-i+1} r_i^{m-i+1}$ , can also be written as the integration of the volume of its  $(m - i)$ -dimensional cross-section—an  $(m - i)$ -ball—that is perpendicular to the  $x_i$  axis and at distance  $x_i$  from the origin as the value  $x_i$  varies from  $-r_i$  to  $r_i$ . Figure 6 gives an illustration. The radius of the  $(m - i)$ -dimensional cross-section is  $r_{i+1} = \sqrt{r_i^2 - x_i^2}$ . Therefore,  $V_{m-i+1} r_i^{m-i+1} = \int_{-r_i}^{r_i} V_{m-i} r_{i+1}^{m-i} dx_i$ , where  $V_{m-i} r_{i+1}^{m-i}$  denotes the volume function of an  $(m - i)$ -ball with radius  $r_{i+1}$ . Inductively, we obtain

$$V_{m-i+1} r_i^{m-i+1} = \int_{-r_i}^{r_i} \int_{-r_{i+1}}^{r_{i+1}} \cdots \int_{-r_m}^{r_m} dx_m \cdots dx_i.$$

We now return to calculating  $\frac{n}{V_m r_0^m} \int_D x_i^4 dx_m \cdots dx_1$ . By symmetry, we can assume that  $i = 1$ .

$$\begin{aligned} \frac{n}{V_m r_0^m} \int_D x_1^4 dx_m \cdots dx_1 &= \frac{n}{V_m r_0^m} \int_{-r_1}^{r_1} \cdots \int_{-r_m}^{r_m} x_1^4 dx_m \cdots dx_1 \\ &= \frac{nV_{m-1}}{V_m r_0^m} \int_{-r_1}^{r_1} x_1^4 r_2^{m-1} dx_1 \\ &= \frac{nV_{m-1}}{V_m r_0^m} \cdot \frac{3\beta(m)r_1^{m+4}}{(m+2)(m+4)} && (\because \text{Lemma 5.3}) \\ &= \frac{3nr_1^{m+4}}{(m+2)(m+4)r_0^m} && (\because (5.4)) \end{aligned}$$

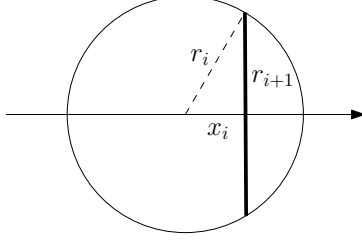


Figure 6: The circle represents an  $(m - i + 1)$ -ball centered at the origin with radius  $r_i$ . The bold segment represents an  $(m - i)$ -dimensional cross-section of the  $(m - i + 1)$ -ball, and its radius  $r_{i+1}$  is equal to  $\sqrt{r_i^2 - x_i^2}$ .

The orthogonal projection of  $\mathcal{M} \cap B_\varepsilon$  onto  $\mathcal{T}$  is contained in the  $m$ -ball centered at the origin with radius  $r_0$ . By (5.8), if we set  $r_1 = r_0$ , then

$$\begin{aligned} \frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^4 \, dV &\leq (1 + O(\varepsilon^2)) \cdot \frac{n}{V_m r_0^m} \int_D x_i^4 \, dx_m \cdots dx_1 \\ &= (1 + O(\varepsilon^2)) \cdot \frac{3nr_0^4}{(m+2)(m+4)} \\ &= \frac{3n\varrho^4 \varepsilon^4}{(m+2)(m+4)} + O(n\varrho^4 \varepsilon^6). \end{aligned}$$

The orthogonal projection of  $\mathcal{M} \cap B_\varepsilon$  onto  $\mathcal{T}$  covers an  $m$ -ball centered at the origin with radius  $(1 - O(\varepsilon^2))r_0$ . Thus, if we set  $r_1$  to be this radius, then

$$\begin{aligned} \frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^4 \, dV &\geq \frac{n}{V_m r_0^m} \int_D x_i^4 \, dx_m \cdots dx_1 \\ &= (1 - O(\varepsilon^2))^{m+4} \cdot \frac{3nr_0^4}{(m+2)(m+4)} \\ &\geq (1 - O(\varepsilon^2)) \cdot \frac{3nr_0^4}{(m+2)(m+4)} \\ &= \frac{3n\varrho^4 \varepsilon^4}{(m+2)(m+4)} - O(n\varrho^4 \varepsilon^6). \end{aligned}$$

Then, it follows from (5.10) that  $\left| \mathbf{E} \left[ \sum_{p=1}^n a_{pi}^4 \right] - \frac{3}{(m+2)(m+4)} n\varrho^4 \varepsilon^4 \right| = O(n\varrho^4 \varepsilon^6)$ , establishing the correctness of (i).

*Analysis of (ii):* The variance of  $\sum_{p=1}^n a_{pi}^2 a_{pj}^2$  equals  $n\mathbf{E}[z_i^4 z_j^4] - n(\mathbf{E}[z_i^2 z_j^2])^2 \leq n\mathbf{E}[z_i^4 z_j^4] \leq n\varrho^8 \varepsilon^8$  because  $|z_i| \leq \varrho\varepsilon$ . Lemma 5.2 implies that  $\left| \sum_{p=1}^n a_{pi}^2 a_{pj}^2 - \mathbf{E} \left[ \sum_{p=1}^n a_{pi}^2 a_{pj}^2 \right] \right| < cn^{2/3} \varrho^4 \varepsilon^4$  with probability  $1 - c^{-2} n^{-1/3}$ . It remains to bound  $\mathbf{E} \left[ \sum_{p=1}^n a_{pi}^2 a_{pj}^2 \right]$ .

Since  $\mathbf{E} \left[ \sum_{p=1}^n a_{pi}^2 a_{pj}^2 \right] = \frac{n}{\text{vol}(\mathcal{M} \cap B_\varepsilon)} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 z_j^2 \, dV$ , we can derive as in (i) the relation  $\frac{(1-O(\varepsilon^2))n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 z_j^2 \, dV \leq \mathbf{E} \left[ \sum_{p=1}^n a_{pi}^2 a_{pj}^2 \right] \leq \frac{(1+O(\varepsilon^2))n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 z_j^2 \, dV$ . Let  $D$  be the  $m$ -ball in  $\mathcal{T}$  centered at the origin with radius  $r_1$ . By symmetry, we can assume that  $i = 1$



and  $j = 2$ .

$$\begin{aligned}
\frac{n}{V_m r_0^m} \int_D x_i^2 x_j^2 dx_m \cdots dx_1 &= \frac{n V_{m-2}}{V_m r_0^m} \int_{-r_1}^{r_1} \int_{-r_2}^{r_2} x_1^2 x_2^2 r_3^{m-2} dx_2 dx_1 \\
&= \frac{n \beta(m-1) \beta(m+2) V_{m-2} r_1^{m+4}}{(m+1)(m+4) V_m r_0^m} \quad (\because \text{Lemma 5.3}) \\
&= \frac{n r_1^{m+4}}{(m+2)(m+4) r_0^m} \quad (\because (5.3) \& (5.5))
\end{aligned}$$

As in (i), if we set  $r_1 = r_0$ , then

$$\begin{aligned}
\frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 z_j^2 dV &\leq (1 + O(\varepsilon^2)) \cdot \frac{n}{V_m r_0^m} \int_D x_i^2 x_j^2 dx_m \cdots dx_1 \\
&= (1 + O(\varepsilon^2)) \cdot \frac{n r_0^4}{(m+2)(m+4)} \\
&= \frac{n \varrho^4 \varepsilon^4}{(m+2)(m+4)} + O(n \varrho^4 \varepsilon^6).
\end{aligned}$$

The orthogonal projection of  $\mathcal{M} \cap B_\varepsilon$  onto  $\mathcal{T}$  covers an  $m$ -ball centered at the origin with radius  $(1 - O(\varepsilon^2))r_0$ . Thus, if we set  $r_1$  to be this radius, then

$$\begin{aligned}
\frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 z_j^2 dV &\geq \frac{n}{V_m r_0^m} \int_D x_i^2 x_j^2 dx_m \cdots dx_1 \\
&= (1 - O(\varepsilon^2))^{m+4} \cdot \frac{n r_0^4}{(m+2)(m+4)} \\
&\geq (1 - O(\varepsilon^2)) \cdot \frac{n r_0^4}{(m+2)(m+4)} \\
&= \frac{n \varrho^4 \varepsilon^4}{(m+2)(m+4)} - O(n \varrho^4 \varepsilon^6).
\end{aligned}$$

Therefore,  $\left| \mathbb{E} \left[ \sum_{p=1}^n a_{pi}^2 a_{pj}^2 \right] - \frac{n \varrho^4 \varepsilon^4}{(m+2)(m+4)} \right| = O(n \varrho^4 \varepsilon^6)$ , establishing the correctness of (ii).

*Analysis of (iii):* The variance of  $\sum_{p=1}^n a_{pi}^2$  is at most  $n \mathbb{E}[z_i^4] \leq n \varrho^4 \varepsilon^4$ . Lemma 5.2 implies that  $\left| \sum_{p=1}^n a_{pi}^2 - \mathbb{E} \left[ \sum_{p=1}^n a_{pi}^2 \right] \right| < c n^{2/3} \varrho^2 \varepsilon^2$  with probability  $1 - c^{-2} n^{-1/3}$ . It remains to bound  $\mathbb{E} \left[ \sum_{p=1}^n a_{pi}^2 \right]$ . Since  $\mathbb{E} \left[ \sum_{p=1}^n a_{pi}^2 \right] = \frac{n}{\text{vol}(\mathcal{M} \cap B_\varepsilon)} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 dV$ , by (5.9),  $\mathbb{E} \left[ \sum_{p=1}^n a_{pi}^2 \right]$  lies between  $(1 - O(\varepsilon^2)) \cdot \frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 dV$  and  $(1 + O(\varepsilon^2)) \cdot \frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 dV$ . Let  $D$  be the  $m$ -ball in  $\mathcal{T}$  centered at the origin with radius  $r_1$ .

$$\begin{aligned}
\frac{n}{V_m r_0^m} \int_D x_i^2 dx_m \cdots dx_1 &= \frac{n}{V_m r_0^m} \int_{-r_1}^{r_1} \cdots \int_{-r_m}^{r_m} x_1^2 dx_m \cdots dx_1 \\
&= \frac{n V_{m-1}}{V_m r_0^m} \int_{-r_1}^{r_1} x_1^2 r_2^{m-1} dx_1 \\
&= \frac{n r_1^{m+2}}{(m+2) r_0^m} \quad (\because \text{Lemma 5.3 and (5.4)})
\end{aligned}$$

If we set  $r_1 = r_0$ , then  $\mathbb{E} \left[ \sum_{p=1}^n a_{pi}^2 \right] \leq (1 + O(\varepsilon^2)) \cdot \frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 dV \leq (1 + O(\varepsilon^2))^2 \cdot \frac{n}{V_m r_0^m} \int_D x_i^2 dx_m \cdots dx_1 = \frac{1+O(\varepsilon^2)}{m+2} \cdot n r_0^2 = \frac{1}{m+2} \cdot n \varrho^2 \varepsilon^2 + O(n \varrho^2 \varepsilon^4)$ . The orthogonal projection

of  $\mathcal{M} \cap B_\varepsilon$  onto  $\mathcal{T}$  covers an  $m$ -ball centered at the origin with radius  $(1 - O(\varepsilon^2))r_0$ . Thus, if we set  $r_1$  to be this radius, then  $\mathbb{E}\left[\sum_{p=1}^n a_{pi}^2\right] \geq (1 - O(\varepsilon^2)) \cdot \frac{n}{V_m r_0^m} \int_{\mathcal{M} \cap B_\varepsilon} z_i^2 dV \geq (1 - O(\varepsilon^2)) \cdot \frac{n}{V_m r_0^m} \int_D x_i^2 dx_m \cdots dx_1 = \frac{(1 - O(\varepsilon^2))^{m+3}}{m+2} \cdot n r_0^2 = \frac{1}{m+2} \cdot n \varrho^2 \varepsilon^2 - O(n \varrho^2 \varepsilon^4)$ . Therefore,  $\left|\mathbb{E}\left[\sum_{p=1}^n a_{pi}^2\right] - \frac{1}{m+2} n \varrho^2 \varepsilon^2\right| = O(n \varrho^2 \varepsilon^4)$ . This establishes the correctness of (iii).

*Analysis of (iv), (v) and (vi):* We prove a more general statement. Consider  $\sum_{p=1}^n \prod_{j=0}^k a_{pi_j}^{e_j}$ , where  $k \geq 1$ ,  $i_j \in [1, m]$  for every  $j \in [0, k]$ ,  $e_0$  is an even (possibly zero) integer, and  $e_j$  is an odd integer for every  $j \in [1, k]$ . We show that  $\left|\sum_{p=1}^n \prod_{j=0}^k a_{pi_j}^{e_j}\right| < cn^{2/3} \varrho^e \varepsilon^e + O(n \varrho^e \varepsilon^{e+2})$  with probability  $1 - c^{-2} n^{-1/3}$ , where  $e = \sum_{j=0}^k e_j$ .

The variance of  $\sum_{p=1}^n \prod_{j=0}^k a_{pi_j}^{e_j}$  is at most  $n \mathbb{E}\left[\prod_{j=0}^k z_{i_j}^{2e_j}\right] \leq n \varrho^{2e} \varepsilon^{2e}$  because  $z_i \leq \varrho \varepsilon$ . By Lemma 5.2, it holds with probability  $1 - c^{-2} n^{-1/3}$  that  $\left|\sum_{p=1}^n \prod_{j=0}^k a_{pi_j}^{e_j}\right| < \left|\mathbb{E}\left[\sum_{p=1}^n \prod_{j=0}^k a_{pi_j}^{e_j}\right]\right| + cn^{2/3} \varrho^e \varepsilon^e$ . Thus, it suffices to bound  $\left|\mathbb{E}\left[\sum_{p=1}^n \prod_{j=0}^k a_{pi_j}^{e_j}\right]\right|$ .

Let  $D$  be the  $m$ -ball in  $\mathcal{T}$  centered at the origin with radius  $r_1 = (1 - O(\varepsilon^2))r_0$  covered by the orthogonal projection of  $\mathcal{M} \cap B_\varepsilon$  onto  $\mathcal{T}$ . Let  $M$  be the portion of  $\mathcal{M} \cap B_\varepsilon$  that projects onto  $D$ . We deal with  $M$  and  $(\mathcal{M} \cap B_\varepsilon) \setminus M$  separately.

Divide  $\mathbb{R}^d$  into  $2^k$  subsets so that for every subset and every  $j \in [1, k]$ , the sign of  $z_{i_j}^{e_j}$  does not flip within the subset. (Note that  $z_{i_0}^{e_0} \geq 0$  as  $e_0$  is even.) Consider two such subsets  $H_k^+ = \{x \in \mathbb{R}^d : x_{i_k}^{e_k} \geq 0 \wedge \forall j \in [1, k-1], x_{i_j}^{e_j} \geq 0\}$  and  $H_k^- = \{x \in \mathbb{R}^d : x_{i_k}^{e_k} \leq 0 \wedge \forall j \in [1, k-1], x_{i_j}^{e_j} \geq 0\}$ . By (5.8),

$$\int_{M \cap H_k^+} \prod_{j=0}^k z_{i_j}^{e_j} dV \leq (1 + O(\varepsilon^2)) \cdot \int_{D \cap H_k^+} \prod_{j=0}^k x_{i_j}^{e_j} dx_m \cdots dx_1.$$

Since  $z_{i_k}^{e_k} \leq 0$  in  $H_k^-$  and  $\int_{M \cap H_k^-} |z_{i_k}^{e_k}| \cdot \prod_{j=0}^{k-1} z_{i_j}^{e_j} dV \geq \int_{D \cap H_k^-} |x_{i_k}^{e_k}| \cdot \prod_{j=0}^{k-1} x_{i_j}^{e_j} dx_m \cdots dx_1$ , we get  $\int_{M \cap H_k^-} \prod_{j=0}^k z_{i_j}^{e_j} dV = - \int_{M \cap H_k^-} |z_{i_k}^{e_k}| \cdot \prod_{j=0}^{k-1} z_{i_j}^{e_j} dV \leq - \int_{D \cap H_k^-} |x_{i_k}^{e_k}| \cdot \prod_{j=0}^{k-1} x_{i_j}^{e_j} dx_m \cdots dx_1$ . By symmetry,  $\int_{D \cap H_k^-} |x_{i_k}^{e_k}| \cdot \prod_{j=0}^{k-1} x_{i_j}^{e_j} dx_m \cdots dx_1 = \int_{D \cap H_k^+} \prod_{j=0}^k x_{i_j}^{e_j} dx_m \cdots dx_1$ . Therefore,

$$\int_{M \cap H_k^-} \prod_{j=0}^k z_{i_j}^{e_j} dV \leq - \int_{D \cap H_k^+} \prod_{j=0}^k x_{i_j}^{e_j} dx_m \cdots dx_1.$$

Let  $H = \{x \in \mathbb{R}^d : \forall j \in [1, k-1], x_{i_j}^{e_j} \geq 0\}$ . It follows that

$$\begin{aligned} \int_{M \cap H} \prod_{j=0}^k z_{i_j}^{e_j} dV &= \int_{M \cap H_k^+} \prod_{j=0}^k z_{i_j}^{e_j} dV + \int_{M \cap H_k^-} \prod_{j=0}^k z_{i_j}^{e_j} dV \\ &\leq O(\varepsilon^2) \cdot \int_{D \cap H_k^+} \prod_{j=0}^k x_{i_j}^{e_j} dx_m \cdots dx_1 \\ &\leq O(\varepsilon^2) \cdot V_m r_1^{m+e} / 2^k \quad (\because x_{i_j} \leq r_1) \\ &\leq O(\varepsilon^2) \cdot V_m r_0^{m+e} / 2^k \quad (\because r_1 \leq r_0) \end{aligned}$$

Conversely,  $\int_{M \cap H_k^+} \prod_{j=0}^k z_{i_j}^{e_j} dV \geq \int_{D \cap H_k^+} \prod_{j=0}^k x_{i_j}^{e_j} dx_m \cdots dx_1$  and  $\int_{M \cap H_k^-} \prod_{j=0}^k z_{i_j}^{e_j} dV = - \int_{M \cap H_k^-} |z_{i_k}^{e_k}| \cdot \prod_{j=0}^{k-1} z_{i_j}^{e_j} dV \geq -(1 + O(\varepsilon^2)) \cdot \int_{D \cap H_k^-} |x_{i_k}^{e_k}| \cdot \prod_{j=0}^{k-1} x_{i_j}^{e_j} dx_m \cdots dx_1$ , which is equal

to  $-(1 + O(\varepsilon^2)) \cdot \int_{D \cap H_k^+} \prod_{j=0}^k x_{i_j}^{e_j} dx_m \cdots dx_1$ , Therefore,

$$\int_{M \cap H} \prod_{j=0}^k z_{i_j}^{e_j} dV \geq -O(\varepsilon^2) \cdot \int_{D \cap H_k^+} \prod_{j=0}^k x_{i_j}^{e_j} dx_m \cdots dx_1 \geq -O(\varepsilon^2) \cdot V_m r_0^{m+e} / 2^k.$$

As a result,  $\left| \int_{M \cap H} \prod_{j=0}^k z_{i_j}^{e_j} dV \right| \leq O(\varepsilon^2) \cdot V_m r_0^{m+e} / 2^k$ .

There are  $2^{k-1} - 1$  other combinations of signs of  $z_{i_j}^{e_j}$  for  $j \in [1, k-1]$ . Each combination gives rise to a subset  $G$  of  $\mathbb{R}^d$  and one can derive as in the above that  $\left| \int_{M \cap G} \prod_{j=0}^k z_{i_j}^{e_j} dV \right| \leq O(\varepsilon^2) \cdot V_m r_0^{m+e} / 2^k$ . Consequently,  $\left| \int_M \prod_{j=0}^k z_{i_j}^{e_j} dV \right| \leq 2^{k-1} \cdot O(\varepsilon^2) \cdot V_m r_0^{m+e} / 2^k = O(\varepsilon^2) \cdot V_m r_0^{m+e}$ .

Let  $D_0$  be the  $m$ -ball in  $\mathcal{T}$  centered at the origin with radius  $r_0$ . Since  $\left| \prod_{j=0}^k z_{i_j}^{e_j} \right| \leq r_0^e$ , by (5.8),  $\left| \int_{(\mathcal{M} \cap B_\varepsilon) \setminus M} \prod_{j=0}^k z_{i_j}^{e_j} dV \right| \leq (1 + O(\varepsilon^2)) \cdot \int_{D_0 \setminus D} r_0^e dx_m \cdots dx_1 = (1 + O(\varepsilon^2)) \cdot V_m r_0^e (r_0^m - r_1^m) = O(\varepsilon^2) \cdot V_m r_0^{m+e}$ .

We conclude that  $\left| \int_{\mathcal{M} \cap B_\varepsilon} \prod_{j=0}^k z_{i_j}^{e_j} dV \right| \leq \left| \int_M \prod_{j=0}^k z_{i_j}^{e_j} dV \right| + \left| \int_{(\mathcal{M} \cap B_\varepsilon) \setminus M} \prod_{j=0}^k z_{i_j}^{e_j} dV \right| = O(\varepsilon^2) \cdot V_m r_0^{m+e}$ . Then,

$$\begin{aligned} \left| \mathbb{E} \left[ \sum_{p=1}^n \prod_{j=0}^k a_{pi_j}^{e_j} \right] \right| &\leq \frac{n}{\text{vol}(\mathcal{M} \cap B_\varepsilon)} \left| \int_{\mathcal{M} \cap B_\varepsilon} \prod_{j=0}^k z_{i_j}^{e_j} dV \right| \\ &\leq \frac{n}{\text{vol}(D)} \cdot O(\varepsilon^2) \cdot V_m r_0^{m+e} \\ &= O(n \varrho^e \varepsilon^{e+2}). \end{aligned}$$

This establishes the correctness of (iv), (v) and (vi).  $\square$

## 6 Eigenvalues of $\mathbf{B}^t \mathbf{B}$

We show in this section that the largest  $m_0$  eigenvalues of  $\mathbf{B}^t \mathbf{B}$  are  $\Theta(n \varrho^4 \varepsilon^4)$  and the largest  $(m_0 + 1)$ -th eigenvalue of  $\mathbf{B}^t \mathbf{B}$  is  $O(n \varrho^4 \varepsilon^6)$ . The bounds on the eigenvalues of  $\mathbf{B}^t \mathbf{B}$  are obtained by proving a series of lemmas using the Gershgorin Circle Theorem [17] and its generalization [13].

### 6.1 Preliminaries

Let  $\mathbf{C}$  be a square matrix. The Gershgorin Circle Theorem states that for each eigenvalue  $\sigma$  of  $\mathbf{C}$ , there exists a row  $i$  of  $\mathbf{C}$  such that  $|\sigma - c_{ii}| \leq \sum_{j \neq i} |c_{ij}|$ . It follows that  $|c_{ii}| - \sum_{j \neq i} |c_{ij}| \leq \sigma \leq \sum_j |c_{ij}|$ . By applying the Gershgorin Circle Theorem to  $\mathbf{C}^t$ , there also exists a column  $j$  of  $\mathbf{C}$  such that  $|\sigma - c_{jj}| \leq \sum_{i \neq j} |c_{ij}|$ , which implies that  $|c_{jj}| - \sum_{i \neq j} |c_{ij}| \leq \sigma \leq \sum_i |c_{ij}|$ . This result has been generalized to the case when  $\mathbf{C}$  is partitioned into blocks [13]. Consider the following partition of  $\mathbf{C}$ :

$$\begin{pmatrix} \mathbf{C}_{11} & \cdots & \mathbf{C}_{1r} \\ \vdots & \ddots & \vdots \\ \mathbf{C}_{r1} & \cdots & \mathbf{C}_{rr} \end{pmatrix} \quad (6.1)$$

That is, there exist integers  $n_i$  such that  $\mathbf{C}_{ij}$  is an  $n_i \times n_j$  matrix. Note that the matrices  $\mathbf{C}_{ii}$ 's are square, but the other  $\mathbf{C}_{ij}$ 's may not be square. Each row of blocks  $(\mathbf{C}_{i1} \cdots \mathbf{C}_{ii} \cdots \mathbf{C}_{ir})$  defines a *generalized gershgorin set*  $G_i$  which contains all real numbers  $\mu$  such that  $\|(\mathbf{C}_{ii} -$

$\mu_{n_i})^{-1}\|^{-1} \leq \sum_{j \neq i} \|C_{ij}\|$ . The eigenvalues of  $C_{ii}$  are in  $G_i$  by a continuity argument [13]. The definition of  $G_i$  implies that:

$$\min\{\mu \in G_i\} \geq \text{smallest eigenvalue of } C_{ii} - \sum_{i \neq j} \|C_{ij}\| \quad (6.2)$$

$$\max\{\mu \in G_i\} \leq \text{largest eigenvalue of } C_{ii} + \sum_{i \neq j} \|C_{ij}\| \quad (6.3)$$

Equations (6.2) and (6.3) help to bound the eigenvalues of  $C$  because  $G_i$  contains some eigenvalues of  $C$  under certain conditions as stated in the following result.

**Lemma 6.1 ([13])** *Consider any partition of a square matrix  $C$  into blocks as in (6.1). Every eigenvalue of  $C$  lies in some generalized gershgorin set  $G_i$  with respect to this partition. Moreover, if a generalized gershgorin set  $G_i$  is disjoint from the union of the other generalized gershgorin sets, then  $G_i$  contains exactly  $n_i$  eigenvalues of  $C$ , where  $n_i$  is the dimension of  $C_{ii}$ .*

In addition to the Gershgorin Circle Theorem, there are also some easy bounds on the 2-norm and eigenvalues of a matrix [17]. For any  $r \times k$  matrix  $U$ ,

$$\|U\| = \|U^t\| \quad \text{and} \quad \|U^t U\| = \|U\|^2.$$

Moreover,

$$\max_{e \in \mathbb{R}^k, \|e\|=1} e^t U^t U e = \sigma_{\max}^2 \quad \text{and} \quad \min_{e \in \mathbb{R}^k, \|e\|=1} e^t U^t U e = \sigma_{\min}^2,$$

where  $\sigma_{\max}$  and  $\sigma_{\min}$  are the largest and smallest singular values of  $U$ . Since  $\|\cdot\|$  satisfies triangle inequality, if  $U = V + W$ , then

$$\|U\| \leq \|V\| + \|W\|.$$

If  $U = (V \ W)$ , where the row dimension of  $U$  is  $r$  and the column dimensions of  $V$  and  $W$  are  $i$  and  $j$ , respectively, then since we can write  $U = (V \ 0_{r,j}) + (0_{r,i} \ W)$ , we get

$$\|U\| \leq \|V\| + \|W\|.$$

If  $U = VW$ , then

$$\|U\| \leq \|V\| \cdot \|W\|.$$

Suppose that the row dimension of  $U$  is  $r$ . Then,  $\|U\| \leq \|U\|_F \leq \sqrt{r} \|U\|$ . Note that  $\|U\|_F = (\sum_{ij} u_{ij}^2)^{1/2} = (\sum_{i=1}^r \|u_{i*}\|^2)^{1/2}$ . Therefore,

$$\|U\| \leq \sqrt{r} \max_{i \in [1,r]} \|u_{i*}\|.$$

Suppose that  $U = V + W$  and all three matrices  $U$ ,  $V$  and  $W$  are symmetric (of dimension  $k$ ) and positive semi-definite. In this case, the minimum eigenvalue of  $U$  is  $\min_{e \in \mathbb{R}^k, \|e\|=1} e^t U e = \min_{e \in \mathbb{R}^k, \|e\|=1} e^t V e + e^t W e$ , which is greater than or equal to both  $\min_{e \in \mathbb{R}^k, \|e\|=1} e^t V e$  and  $\min_{e \in \mathbb{R}^k, \|e\|=1} e^t W e$ . Therefore, if  $U = V + W$  for some symmetric and positive semi-definite matrices  $U$ ,  $V$  and  $W$ , then

$$\text{min. eigenvalue of } U \geq \max\{\text{min. eigenvalue of } V, \text{min. eigenvalue of } W\}.$$

To facilitate the analysis of the eigenvalues of  $B^t B$ , it is convenient to assume that the coordinate axes  $x_1, \dots, x_m$  span the tangent space  $\mathcal{T}$  of  $\mathcal{M}$  at the origin. That is, for every  $p \in [1, n]$  and every  $i \in [1, m]$ ,  $a_{pi}$  is the coordinate of  $a_p$  on the  $x_i$ -axis. The following result shows that the eigenvalues of  $B^t B$  are preserved by rotations in  $\mathbb{R}^d$  that keep the origin fixed.

**Lemma 6.2** *The eigenvalues of  $\mathbf{B}^t\mathbf{B}$  are preserved by rotations in  $\mathbb{R}^d$  that keep the origin fixed.*

*Proof.* Recall that  $d_0 = \binom{d+1}{2}$ . Each row of  $\mathbf{B}$  is a vector in  $\mathbb{R}^{d_0}$  and it is the image of the function  $h : (y_1 \ \cdots \ y_d)^t \mapsto \left( \frac{1}{\sqrt{2}}y_1^2 \ y_1y_2 \ \cdots \ y_1y_d \ \frac{1}{\sqrt{2}}y_2^2 \ y_2y_3 \ \cdots \ \frac{1}{\sqrt{2}}y_d^2 \right)^t$ . The linear space spanned by  $h(\mathbb{R}^d)$  is  $\mathbb{R}^{d_0}$ . To see this, consider the vectors  $\mathbf{e}_{ij}$ , where  $i \in [1, d]$  and  $j \in [i, d]$ , such that the  $i$ -th and  $j$ -th entries of  $\mathbf{e}_{ij}$  are ones and all other entries of  $\mathbf{e}_{ij}$  are zeros. There are  $d_0$  such  $\mathbf{e}_{ij}$ 's and the vectors  $h(\mathbf{e}_{ij})$ 's are linearly independent because for any  $k < l$ ,  $h(\mathbf{e}_{kl})$  contains a 1 in the position of  $y_ky_l$  and no other  $h(\mathbf{e}_{ij})$ 's do.

Take any rotation  $R$  in  $\mathbb{R}^d$  that keeps the origin fixed. Define the transformation  $\varphi : h(\mathbf{v}) \mapsto h \circ R(\mathbf{v})$ . For every pair of vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ ,  $h(\mathbf{u})^t \cdot h(\mathbf{v}) = \frac{1}{2} \sum_{i=1}^d u_i^2 v_i^2 + \sum_{i=1}^d \sum_{j=i+1}^d u_i u_j v_i v_j = \frac{1}{2} (\sum_{i=1}^d u_i v_i)^2 = \frac{1}{2} (\mathbf{u}^t \cdot \mathbf{v})^2$ . Since  $R$  preserves distances and angles in  $\mathbb{R}^d$ ,  $\mathbf{u}^t \cdot \mathbf{v} = R(\mathbf{u})^t \cdot R(\mathbf{v})$ , which implies that  $(\varphi \circ h(\mathbf{u}))^t \cdot \varphi \circ h(\mathbf{v}) = h(\mathbf{u})^t \cdot h(\mathbf{v})$ . That is,  $\varphi$  preserves distances and angles in  $h(\mathbb{R}^d)$ .

Since the  $h(\mathbf{e}_{ij})$ 's form a basis of  $\mathbb{R}^{d_0}$ , we can define a linear transformation  $\psi$  in  $\mathbb{R}^{d_0}$  such that  $\psi \circ h(\mathbf{e}_{ij}) = \varphi \circ h(\mathbf{e}_{ij})$  for every  $i \in [1, d]$  and every  $j \in [i, d]$ . For any  $i, k \in [1, d]$ , any  $j \in [i, d]$  and any  $l \in [k, d]$ ,  $(\psi \circ h(\mathbf{e}_{ij}))^t \cdot \psi \circ h(\mathbf{e}_{kl}) = (\varphi \circ h(\mathbf{e}_{ij}))^t \cdot \varphi \circ h(\mathbf{e}_{kl}) = h(\mathbf{e}_{ij})^t \cdot h(\mathbf{e}_{kl})$ , which implies that  $\psi$  preserves distances and angles in  $\mathbb{R}^{d_0}$  and hence  $\psi$  is an isometry in  $\mathbb{R}^{d_0}$ . Since  $\psi$  and  $\varphi$  agree on the  $h(\mathbf{e}_{ij})$ 's by definition and both  $\psi$  and  $\varphi$  preserve distances in  $h(\mathbb{R}^d)$ ,  $\psi(\mathbf{z})$  must be equal to  $\varphi(\mathbf{z})$  for every vector  $\mathbf{z} \in h(\mathbb{R}^d)$ .

We conclude that the effect on  $h(\mathbb{R}^d)$  caused by the rotation  $R$  in  $\mathbb{R}^d$  is produced by the isometry  $\psi$  in  $\mathbb{R}^{d_0}$ . Since the eigenvalues of  $\mathbf{B}^t\mathbf{B}$  are invariant under isometries in  $\mathbb{R}^{d_0}$ , they are not changed by the rotation  $R$ .  $\square$

## 6.2 Analysis

Suppose that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . There are two kinds of columns in  $\mathbf{B}$ , namely, the ‘‘double’’ columns  $\left( \frac{1}{\sqrt{2}}a_{1i}^2 \ \cdots \ \frac{1}{\sqrt{2}}a_{ni}^2 \right)^t$  for  $i \in [1, d]$  and the ‘‘cross’’ columns  $(a_{1i}a_{1j} \ \cdots \ a_{ni}a_{nj})^t$  for  $i \in [1, d]$  and  $j \in [i+1, d]$ . Rearranging the columns in  $\mathbf{B}$  does not change the eigenvalues of  $\mathbf{B}^t\mathbf{B}$ . For convenience, we rearrange the columns of  $\mathbf{B}$  so that  $\mathbf{B} = (\mathbf{B}_{TT} \ \mathbf{B}_{TN} \ \mathbf{B}_{NN})$ , where  $\mathbf{B}_{TT}$  consists of the ‘‘double’’ columns for  $i \in [1, m]$  and the ‘‘cross’’ columns for  $i \in [1, m]$  and  $j \in [i+1, m]$ ,  $\mathbf{B}_{TN}$  consists of the ‘‘cross’’ columns for  $i \in [1, m]$  and  $j \in [m+1, d]$ , and  $\mathbf{B}_{NN}$  consists of the ‘‘double’’ columns for  $i \in [m+1, d]$  and the ‘‘cross’’ columns for  $i \in [m+1, d]$  and  $j \in [i+1, d]$ .

Recall that  $d_0 = \binom{d+1}{2}$  and  $m_0 = \binom{m+1}{2}$ .  $\mathbf{B}_{TT}$  has  $m_0$  columns,  $\mathbf{B}_{TN}$  has  $m(d-m)$  columns, and  $\mathbf{B}_{NN}$  contains  $d_0 - m_0 - md + m^2$  columns. The matrix  $\mathbf{B}^t\mathbf{B}$  can be divided into blocks as follows.

$$\mathbf{B}^t\mathbf{B} = \begin{pmatrix} \mathbf{B}_{TT}^t\mathbf{B}_{TT} & \mathbf{B}_{TT}^t\mathbf{B}_{TN} & \mathbf{B}_{TT}^t\mathbf{B}_{NN} \\ \mathbf{B}_{TN}^t\mathbf{B}_{TT} & \mathbf{B}_{TN}^t\mathbf{B}_{TN} & \mathbf{B}_{TN}^t\mathbf{B}_{NN} \\ \mathbf{B}_{NN}^t\mathbf{B}_{TT} & \mathbf{B}_{NN}^t\mathbf{B}_{TN} & \mathbf{B}_{NN}^t\mathbf{B}_{NN} \end{pmatrix} \quad (6.4)$$

We first analyze the eigenvalues of  $\mathbf{B}_{TT}^t\mathbf{B}_{TT}$  and the singular values of  $\mathbf{B}_{TN}$  and  $\mathbf{B}_{NN}$ .

**Lemma 6.3** *Assume that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , the eigenvalues of  $\mathbf{B}_{TT}^t\mathbf{B}_{TT}$  are  $\Theta(n\rho^4\varepsilon^4)$  and so  $\|\mathbf{B}_{TT}\| = \Theta(\sqrt{n}\rho^2\varepsilon^2)$ .*

*Proof.* For every  $p \in [1, n]$  and every  $i \in [1, m]$ ,  $|a_{pi}| \leq \rho\varepsilon$ . Thus, for every  $p \in [1, n]$  and every quadruple of possibly non-distinct  $i, j, k, l \in [1, m]$ ,  $\sum_{p=1}^n |a_{pi}a_{pj}a_{pk}a_{pl}| \leq n\rho^4\varepsilon^4$ . It follows that

the maximum absolute row sum of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  is at most  $m_0 n \varrho^4 \varepsilon^4 = O(n \varrho^4 \varepsilon^4)$ , which is an upper bound on the largest eigenvalue by the Gershgorin Circle Theorem.

To prove the lower bound, rearrange the columns of  $\mathbf{B}_{TT}$  such that its leftmost  $n \times m$  submatrix is  $\begin{pmatrix} \frac{1}{\sqrt{2}} a_{11}^2 & \cdots & \frac{1}{\sqrt{2}} a_{1m}^2 \\ \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{2}} a_{n1}^2 & \cdots & \frac{1}{\sqrt{2}} a_{nm}^2 \end{pmatrix}$ . This rearrangement does not change the eigenvalues of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$ . Let  $\mathbf{V}$  be the trailing  $(m_0 - m) \times (m_0 - m)$  submatrix of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$ . Then,

$$\mathbf{B}_{TT}^t \mathbf{B}_{TT} = \left( \begin{array}{cccc|c} \frac{1}{2} \sum_{p=1}^n a_{p1}^4 & \frac{1}{2} \sum_{p=1}^n a_{p1}^2 a_{p2}^2 & \cdots & \frac{1}{2} \sum_{p=1}^n a_{p1}^2 a_{pm}^2 & * \\ \frac{1}{2} \sum_{p=1}^n a_{p1}^2 a_{p2}^2 & \frac{1}{2} \sum_{p=1}^n a_{p2}^4 & \cdots & \frac{1}{2} \sum_{p=1}^n a_{p2}^2 a_{pm}^2 & * \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \frac{1}{2} \sum_{p=1}^n a_{p1}^2 a_{pm}^2 & \frac{1}{2} \sum_{p=1}^n a_{p2}^2 a_{pm}^2 & \cdots & \frac{1}{2} \sum_{p=1}^n a_{pm}^4 & * \\ * & * & \cdots & * & \mathbf{V} \end{array} \right) \quad (6.5)$$

Define an  $m \times m$  matrix  $\mathbf{W}$  whose entries are identical and equal to  $\frac{1}{2(m+2)(m+4)} n \varrho^4 \varepsilon^4$ . Define another  $m \times m$  matrix  $\mathbf{U}$  such that  $u_{ij}$  is equal to the  $(i, j)$  entry of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  minus  $\frac{1}{2(m+2)(m+4)} n \varrho^4 \varepsilon^4$ . We split  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  into the following sum of matrices.

$$\mathbf{B}_{TT}^t \mathbf{B}_{TT} = \begin{pmatrix} \mathbf{U} & * \\ * & \mathbf{V} \end{pmatrix} + \begin{pmatrix} \mathbf{W} & \mathbf{0}_{m, m_0-m} \\ \mathbf{0}_{m_0-m, m} & \mathbf{0}_{m_0-m, m_0-m} \end{pmatrix} \quad (6.6)$$

All matrices in (6.6) are symmetric, and  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  is clearly positive semi-definite. We show that the two matrices on the right hand side of (6.6) are also positive semi-definite by bounding their eigenvalues from below. Then, we can conclude that the minimum eigenvalue of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  is at least the maximum of the minimum eigenvalues of the two matrices on the right hand side of (6.6). Since the entries of  $\mathbf{W}$  are identical, the matrix  $\begin{pmatrix} \mathbf{W} & \mathbf{0}_{m, m_0-m} \\ \mathbf{0}_{m_0-m, m} & \mathbf{0}_{m_0-m, m_0-m} \end{pmatrix}$  has rank one.

One can verify that  $(1/\sqrt{m} \ \cdots \ 1/\sqrt{m} \ \mathbf{0}_{m_0-m})^t$  is a unit eigenvector and the only non-zero eigenvalue is  $\frac{m}{2(m+2)(m+4)} n \varrho^4 \varepsilon^4$ . So the minimum eigenvalue of  $\begin{pmatrix} \mathbf{W} & \mathbf{0}_{m, m_0-m} \\ \mathbf{0}_{m_0-m, m} & \mathbf{0}_{m_0-m, m_0-m} \end{pmatrix}$  is zero. It remains to bound the minimum eigenvalue of  $\begin{pmatrix} \mathbf{U} & * \\ * & \mathbf{V} \end{pmatrix}$  from below.

Apply Lemma 5.4(i) with the constant  $c = \frac{2}{3(m+2)(m+4)}$ . Then, it follows from (6.5) that, with probability  $1 - O(n^{-1/3})$ , for every  $i \in [1, m]$ , the  $(i, i)$  entry of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  is at least  $\frac{3}{2(m+2)(m+4)} n \varrho^4 \varepsilon^4 - \frac{1}{3(m+2)(m+4)} n^{2/3} \varrho^4 \varepsilon^4 - O(n \varrho^4 \varepsilon^6)$ , which implies that

$$\begin{aligned} |u_{ii}| &\geq \frac{3n\varrho^4\varepsilon^4}{2(m+2)(m+4)} - \frac{n\varrho^4\varepsilon^4}{2(m+2)(m+4)} - \frac{n^{2/3}\varrho^4\varepsilon^4}{3(m+2)(m+4)} - O(n\varrho^4\varepsilon^6) \\ &= \frac{n\varrho^4\varepsilon^4}{(m+2)(m+4)} - \frac{n^{2/3}\varrho^4\varepsilon^4}{3(m+2)(m+4)} - O(n\varrho^4\varepsilon^6). \end{aligned} \quad (6.7)$$

Apply Lemma 5.4(ii) with the constant  $c = \frac{2}{3m_0(m+2)(m+4)}$ . Then, it follows from (6.5) that, with probability  $1 - O(n^{-1/3})$ , for every pair of distinct  $i, j \in [1, m]$ , the  $(i, j)$  entry of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  is at most  $\frac{1}{2(m+2)(m+4)} n \varrho^4 \varepsilon^4 + \frac{1}{3m_0(m+2)(m+4)} n^{2/3} \varrho^4 \varepsilon^4 + O(n \varrho^4 \varepsilon^6)$ , which implies that

$$\begin{aligned} |u_{ij}| &\leq \frac{n\varrho^4\varepsilon^4}{2(m+2)(m+4)} - \frac{n\varrho^4\varepsilon^4}{2(m+2)(m+4)} + \frac{n^{2/3}\varrho^4\varepsilon^4}{3m_0(m+2)(m+4)} + O(n\varrho^4\varepsilon^6) \\ &= \frac{n^{2/3}\varrho^4\varepsilon^4}{3m_0(m+2)(m+4)} + O(n\varrho^4\varepsilon^6). \end{aligned} \quad (6.8)$$

By the Gershgorin Circle Theorem, every eigenvalue of  $\begin{pmatrix} \mathbf{U} & * \\ * & \mathbf{V} \end{pmatrix}$  is at least the absolute value of the diagonal entry minus the absolute values of the off-diagonal entries in the same row for some row.

Consider a row that contains the entry  $u_{ii}$  of  $\mathbf{U}$  for some  $i \in [1, m]$ . By (6.8), the absolute value of each off-diagonal entry in  $\mathbf{U}$  is at most  $\frac{1}{3m_0(m+2)(m+4)}n^{2/3}\varrho^4\varepsilon^4 + O(n\varrho^4\varepsilon^6)$ . The other off-diagonal entries are  $\frac{1}{\sqrt{2}}\sum_{p=1}^n a_{pi}^3 a_{pj}$  and  $\frac{1}{\sqrt{2}}\sum_{p=1}^n a_{pi}^2 a_{pj} a_{pk}$  for some distinct  $i, j, k \in [1, m]$ . Apply Lemma 5.4(iv) with the constant  $c = \frac{\sqrt{2}}{3m_0(m+2)(m+4)}$ . It implies that  $\left| \frac{1}{\sqrt{2}}\sum_{p=1}^n a_{pi}^3 a_{pj} \right|$  and  $\left| \frac{1}{\sqrt{2}}\sum_{p=1}^n a_{pi}^2 a_{pj} a_{pk} \right|$  are at most  $\frac{1}{3m_0(m+2)(m+4)}n^{2/3}\varrho^4\varepsilon^4 + O(n\varrho^4\varepsilon^6)$  with probability  $1 - O(n^{-1/3})$ . By (6.7),  $|u_{ii}|$  minus the absolute values of the off-diagonal entries in the  $i$ th row is greater than  $\frac{1}{(m+2)(m+4)}n\varrho^4\varepsilon^4 - \frac{2}{3(m+2)(m+4)}n^{2/3}\varrho^4\varepsilon^4 - O(n\varrho^4\varepsilon^6)$ , which is  $\Omega(n\varrho^4\varepsilon^4)$ .

Consider a row that contains a diagonal entry of  $\mathbf{V}$ . This diagonal entry of  $\mathbf{V}$  is  $\sum_{p=1}^n a_{pi}^2 a_{pj}^2$  for some distinct  $i, j \in [1, m]$ . Apply Lemma 5.4(ii) with the constant  $c = \frac{1}{3(m+2)(m+4)}$ . It implies that  $\sum_{p=1}^n a_{pi}^2 a_{pj}^2 \geq \frac{1}{(m+2)(m+4)}n\varrho^4\varepsilon^4 - \frac{1}{3(m+2)(m+4)}n^{2/3}\varrho^4\varepsilon^4 - O(n\varrho^4\varepsilon^6)$  with probability  $1 - O(n^{-1/3})$ . The off-diagonal entries are  $\frac{1}{\sqrt{2}}\sum_{p=1}^n a_{pi}^3 a_{pj}$ ,  $\frac{1}{\sqrt{2}}\sum_{p=1}^n a_{pi}^2 a_{pj} a_{pk}$  and  $\sum_{p=1}^n a_{pi} a_{pj} a_{pk} a_{pl}$  for some distinct  $i, j, k, l \in [1, m]$ . A similar analysis as in the previous paragraph shows that  $\sum_{p=1}^n a_{pi}^2 a_{pj}^2$  minus the off-diagonal entries in the same row is  $\Omega(n\varrho^4\varepsilon^4)$ .  $\square$

**Lemma 6.4** *Assume that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ .  $\|\mathbf{B}_{TN}\| = O(\sqrt{n}\varrho^2\varepsilon^3)$  and  $\|\mathbf{B}_{NN}\| = O(\sqrt{n}\varrho^2\varepsilon^4)$ .*

*Proof.* For every  $p \in [1, n]$ ,  $(\sum_{i=1}^m a_{pi}^2)^{1/2} \leq \|\mathbf{a}_p\| \leq \varrho\varepsilon$ . The 2-norm of a row of  $\mathbf{B}_{TN}$  is  $(\sum_{i=1}^m a_{pi}^2 \sum_{j=m+1}^d a_{pj}^2)^{1/2} = (\sum_{i=1}^m a_{pi}^2)^{1/2} \cdot (\sum_{j=m+1}^d a_{pj}^2)^{1/2}$ . The distance from  $\mathbf{a}_p$  to  $\mathcal{T}$  is known to be  $O(\varrho\varepsilon^2)$  [15, Lemma 6], so  $(\sum_{j=m+1}^d a_{pj}^2)^{1/2} = O(\varrho\varepsilon^2)$ . It follows that the 2-norm of a row of  $\mathbf{B}_{TN}$  is  $O(\varrho^2\varepsilon^3)$ . There are  $n$  rows in  $\mathbf{B}_{TN}$ , implying that  $\|\mathbf{B}_{TN}\| = O(\sqrt{n}\varrho^2\varepsilon^3)$ . The 2-norm of a row of  $\mathbf{B}_{NN}$  is no more than  $(\sum_{i=m+1}^d a_{pi}^2 \sum_{j=i}^d a_{pj}^2)^{1/2} \leq \sum_{j=m+1}^d a_{pj}^2 = O(\varrho^2\varepsilon^4)$ . Summing over the  $n$  rows in  $\mathbf{B}_{NN}$  shows that  $\|\mathbf{B}_{NN}\| = O(\sqrt{n}\varrho^2\varepsilon^4)$ .  $\square$

We are now ready to give bounds on the eigenvalues of  $\mathbf{B}^t\mathbf{B}$ . The bound on the  $d_0 - m_0$  smallest eigenvalues is not the best possible yet. We bootstrap a better bound from it later.

**Lemma 6.5** *If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , the  $m_0$  largest eigenvalues of  $\mathbf{B}^t\mathbf{B}$  are  $\Theta(n\varrho^4\varepsilon^4)$  and the  $d_0 - m_0$  smallest eigenvalues of  $\mathbf{B}^t\mathbf{B}$  are  $O(n\varrho^4\varepsilon^5)$ .*

*Proof.* By Lemma 6.2, we can rotate  $\mathbb{R}^d$  so that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$  because the eigenvalues of  $\mathbf{B}^t\mathbf{B}$  are not affected. Then, we can partition  $\mathbf{B}^t\mathbf{B}$  as shown in (6.4). We define three generalized gershgorin sets [13] as follows.

The first set  $G_1$  is for the row of blocks  $(\mathbf{B}_{TT}^t \mathbf{B}_{TT} \quad \mathbf{B}_{TT}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TT}^t \mathbf{B}_{NN})$ . By (6.2) and (6.3), the real numbers in  $G_1$  are at least the minimum eigenvalue of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  minus  $\|\mathbf{B}_{TT}^t \mathbf{B}_{TN}\| + \|\mathbf{B}_{TT}^t \mathbf{B}_{NN}\|$  and at most the maximum eigenvalue of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  plus  $\|\mathbf{B}_{TT}^t \mathbf{B}_{TN}\| + \|\mathbf{B}_{TT}^t \mathbf{B}_{NN}\|$ . By Lemma 6.3, with probability  $1 - O(n^{-1/3})$ ,  $\|\mathbf{B}_{TT}^t \mathbf{B}_{TT}\| = \Theta(n\varrho^4\varepsilon^4)$  and  $\|\mathbf{B}_{TT}^t\| = \Theta(\sqrt{n}\varrho^2\varepsilon^2)$ . Then, it follows from Lemma 6.4 that, with probability  $1 - O(n^{-1/3})$ ,  $\|\mathbf{B}_{TT}^t \mathbf{B}_{TN}\| \leq \|\mathbf{B}_{TT}^t\| \cdot \|\mathbf{B}_{TN}\| = O(n\varrho^4\varepsilon^5)$  and  $\|\mathbf{B}_{TT}^t \mathbf{B}_{NN}\| \leq \|\mathbf{B}_{TT}^t\| \cdot \|\mathbf{B}_{NN}\| = O(n\varrho^4\varepsilon^6)$ . Thus, the numbers in  $G_1$  are  $\Theta(n\varrho^4\varepsilon^4)$  with probability  $1 - O(n^{-1/3})$ .

The second set  $G_2$  is for the row of blocks  $(\mathbf{B}_{TN}^t \mathbf{B}_{TT} \quad \mathbf{B}_{TN}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TN}^t \mathbf{B}_{NN})$ . By (6.3), the real numbers in  $G_2$  are at most the maximum eigenvalue of  $\mathbf{B}_{TN}^t \mathbf{B}_{TN}$  plus  $\|\mathbf{B}_{TN}^t \mathbf{B}_{TT}\| +$

$\|\mathbf{B}_{TN}^t \mathbf{B}_{NN}\|$ . By Lemmas 6.3 and 6.4, with probability  $1 - O(n^{-1/3})$ ,  $\|\mathbf{B}_{TN}^t \mathbf{B}_{TN}\| = \|\mathbf{B}_{TN}\|^2 = O(n\varrho^4 \varepsilon^6)$  and  $\|\mathbf{B}_{TN}^t \mathbf{B}_{TT}\| + \|\mathbf{B}_{TN}^t \mathbf{B}_{NN}\| \leq \|\mathbf{B}_{TN}\| \cdot \|\mathbf{B}_{TT}\| + \|\mathbf{B}_{TN}\| \cdot \|\mathbf{B}_{NN}\| = O(n\varrho^4 \varepsilon^5)$ . Thus, the numbers in  $G_2$  are  $O(n\varrho^4 \varepsilon^5)$  with probability  $1 - O(n^{-1/3})$ .

The third set  $G_3$  is for the row of blocks  $(\mathbf{B}_{NN}^t \mathbf{B}_{TT} \quad \mathbf{B}_{NN}^t \mathbf{B}_{TN} \quad \mathbf{B}_{NN}^t \mathbf{B}_{NN})$ . By (6.3), the real numbers in  $G_3$  are at most the maximum eigenvalue of  $\mathbf{B}_{NN}^t \mathbf{B}_{NN}$  plus  $\|\mathbf{B}_{NN}^t \mathbf{B}_{TT}\| + \|\mathbf{B}_{NN}^t \mathbf{B}_{TN}\|$ . By Lemmas 6.3 and 6.4, with probability  $1 - O(n^{-1/3})$ ,  $\|\mathbf{B}_{NN}^t \mathbf{B}_{NN}\| = \|\mathbf{B}_{NN}\|^2 = O(n\varrho^4 \varepsilon^8)$  and  $\|\mathbf{B}_{NN}^t \mathbf{B}_{TT}\| + \|\mathbf{B}_{NN}^t \mathbf{B}_{TN}\| \leq \|\mathbf{B}_{NN}\| \cdot \|\mathbf{B}_{TT}\| + \|\mathbf{B}_{NN}\| \cdot \|\mathbf{B}_{TN}\| = O(n\varrho^4 \varepsilon^6)$ . Thus, the numbers in  $G_3$  are  $O(n\varrho^4 \varepsilon^6)$  with probability  $1 - O(n^{-1/3})$ .

If  $\varepsilon$  is sufficiently small, the numbers in  $G_1$  are much bigger than those in  $G_2$  and  $G_3$ , implying that  $G_1 \cap (G_2 \cup G_3)$  is empty. By Lemma 6.1, the disjointness of  $G_1$  from  $G_2 \cup G_3$  implies that  $G_1$  and  $G_2 \cup G_3$  contain exactly  $m_0$  and  $d_0 - m_0$  eigenvalues of  $\mathbf{B}^t \mathbf{B}$ , respectively. Hence,  $G_1$  contains the  $m_0$  largest eigenvalues and  $G_2 \cup G_3$  contains the  $d_0 - m_0$  smallest eigenvalues.  $\square$

Lemma 6.5 allows us to show a tighter bound  $O(n\varrho^4 \varepsilon^6)$  on the  $(m_0 + 1)$ -th largest eigenvalue of  $\mathbf{B}^t \mathbf{B}$ . It will be important later that this bound is smaller than the bound on the  $m_0$  largest eigenvalues by a factor  $\varepsilon^2$ .

**Lemma 6.6** *If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , the  $(m_0 + 1)$ -th largest eigenvalue of  $\mathbf{B}^t \mathbf{B}$  is  $O(n\varrho^4 \varepsilon^6)$ .*

*Proof.* By Lemma 6.2, we can rotate  $\mathbb{R}^d$  so that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$  and then partition  $\mathbf{B}^t \mathbf{B}$  as shown in (6.4). Let  $\sigma$  be an eigenvalue of  $\mathbf{B}^t \mathbf{B}$  among the  $d_0 - m_0$  smallest ones. By Lemma 6.5,  $\sigma$  is  $O(n\varrho^4 \varepsilon^5)$  with probability  $1 - O(n^{-1/3})$ . Let  $\mathbf{e}$  be an eigenvector of  $\mathbf{B}^t \mathbf{B}$  corresponding to  $\sigma$ . Divide  $\mathbf{e}$  into two parts  $(\mathbf{v}^t \quad \mathbf{w}^t)^t$ , where  $\mathbf{v}$  consists of the first  $m_0$  coordinates of  $\mathbf{e}$  and  $\mathbf{w}$  consists of the last  $d_0 - m_0$  coordinates of  $\mathbf{e}$ .

We claim that  $\mathbf{w} \neq \mathbf{0}_{d_0 - m_0, 1}$  with probability  $1 - O(n^{-1/3})$ . If  $\mathbf{w} = \mathbf{0}_{d_0 - m_0, 1}$ , then the following relation holds as  $\mathbf{e}$  is an eigenvector of  $\mathbf{B}^t \mathbf{B}$ .

$$\mathbf{B}^t \mathbf{B} \cdot \mathbf{e} = \begin{pmatrix} \mathbf{B}_{TT}^t \mathbf{B}_{TT} \cdot \mathbf{v} \\ * \end{pmatrix} = \sigma \mathbf{e} = \begin{pmatrix} \sigma \mathbf{v} \\ \mathbf{0}_{d_0 - m_0, 1} \end{pmatrix}.$$

This implies that  $\sigma$  is an eigenvalue of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$ . Then, either  $\sigma$  is not  $O(n\varrho^4 \varepsilon^5)$  which occurs with probability  $O(n^{-1/3})$  by Lemma 6.5, or  $\sigma = O(n\varrho^4 \varepsilon^5)$  is an eigenvalue of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  which occurs with probability  $O(n^{-1/3})$  by Lemma 6.3. This proves our claim.

From now on, assume that  $\mathbf{w} \neq \mathbf{0}_{d_0 - m_0, 1}$  and  $\mathbf{e}$  is scaled such that  $\|\mathbf{w}\| = 1$ .

Next, we show that  $\|\mathbf{v}\| = O(\varepsilon)$ . Refer to the partition of  $\mathbf{B}^t \mathbf{B}$  in (6.4). Expanding the equation  $\mathbf{B}^t \mathbf{B} \cdot \mathbf{e} = \sigma \mathbf{e}$  gives the equation  $\mathbf{B}_{TT}^t \mathbf{B}_{TT} \cdot \mathbf{v} + (\mathbf{B}_{TT}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TT}^t \mathbf{B}_{NN}) \cdot \mathbf{w} = \sigma \mathbf{v}$ . It implies that  $\mathbf{v} = -(\mathbf{B}_{TT}^t \mathbf{B}_{TT} - \sigma \mathbf{I}_{m_0})^{-1} \cdot (\mathbf{B}_{TT}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TT}^t \mathbf{B}_{NN}) \cdot \mathbf{w}$ . Therefore,

$$\|\mathbf{v}\| \leq \|(\mathbf{B}_{TT}^t \mathbf{B}_{TT} - \sigma \mathbf{I}_{m_0})^{-1}\| \cdot \|(\mathbf{B}_{TT}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TT}^t \mathbf{B}_{NN})\|. \quad (6.9)$$

Note that  $\|(\mathbf{B}_{TT}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TT}^t \mathbf{B}_{NN})\| \leq \|\mathbf{B}_{TT}^t \mathbf{B}_{TN}\| + \|\mathbf{B}_{TT}^t \mathbf{B}_{NN}\| \leq \|\mathbf{B}_{TT}\| \cdot \|\mathbf{B}_{TN}\| + \|\mathbf{B}_{TT}\| \cdot \|\mathbf{B}_{NN}\|$ . Then, Lemmas 6.3 and 6.4 imply that  $\|(\mathbf{B}_{TT}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TT}^t \mathbf{B}_{NN})\| = O(n\varrho^4 \varepsilon^5)$  with probability  $1 - O(n^{-1/3})$ . By Lemmas 6.3 and 6.5, with probability  $1 - O(n^{-1/3})$ , the eigenvalues of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT} - \sigma \mathbf{I}_{m_0}$  are  $\Theta(n\varrho^4 \varepsilon^4) - O(n\varrho^4 \varepsilon^5) = \Theta(n\varrho^4 \varepsilon^4)$ , implying that  $\|(\mathbf{B}_{TT}^t \mathbf{B}_{TT} - \sigma \mathbf{I}_{m_0})^{-1}\| = \Theta(1/(n\varrho^4 \varepsilon^4))$ . Plugging the bounds on  $\|(\mathbf{B}_{TT}^t \mathbf{B}_{TN} \quad \mathbf{B}_{TT}^t \mathbf{B}_{NN})\|$  and  $\|(\mathbf{B}_{TT}^t \mathbf{B}_{TT} - \sigma \mathbf{I}_{m_0})^{-1}\|$  into (6.9) gives  $\|\mathbf{v}\| = O(\varepsilon)$ .

By the definition of eigenvalues, if we project each row vector of  $\mathbf{B}$  onto the support line of  $\mathbf{e}$ , the sum of the squared lengths of the projections is  $\sigma$ . We show that this sum is  $O(n\varrho^4 \varepsilon^6)$  as follows.



Take the  $p$ th row  $\mathbf{b}_{p^*}$  of  $\mathbf{B}$ . Divide  $\mathbf{b}_{p^*}$  into two parts  $(\bar{\mathbf{b}}_{p^*}, \tilde{\mathbf{b}}_{p^*})$ , where  $\bar{\mathbf{b}}_{p^*}$  consists of the first  $m_0$  entries and  $\tilde{\mathbf{b}}_{p^*}$  consists of the last  $d_0 - m_0$  entries. Note that  $\|\bar{\mathbf{b}}_{p^*}\| = \sqrt{\frac{1}{2} \sum_{i=1}^m a_{pi}^4 + \sum_{i=1}^m \sum_{j=i+1}^m a_i^2 a_j^2} = \sqrt{\frac{1}{2} \left( \sum_{i=1}^m a_{pi}^2 \right)^2} < \|\mathbf{a}_p\|^2 \leq \varrho^2 \varepsilon^2$ . Therefore,

$$\mathbf{v}^t \cdot \bar{\mathbf{b}}_{p^*} \leq \|\mathbf{v}\| \cdot \|\bar{\mathbf{b}}_{p^*}\| \leq O(\varepsilon) \cdot \varrho^2 \varepsilon^2 = O(\varrho^2 \varepsilon^3).$$

By grouping terms in  $\|\tilde{\mathbf{b}}_{p^*}\|^2$ , we get

$$\begin{aligned} \|\tilde{\mathbf{b}}_{p^*}\|^2 &= \sum_{i=1}^m \sum_{j=m+1}^d a_{pi}^2 a_{pj}^2 + \sum_{i=m+1}^d \frac{1}{2} a_{pi}^4 + \sum_{i=m+1}^d \sum_{j=i+1}^d a_{pi}^2 a_{pj}^2 \\ &= \left( \sum_{i=1}^m a_{pi}^2 \right) \left( \sum_{j=m+1}^d a_{pj}^2 \right) + \frac{1}{2} \left( \sum_{j=m+1}^d a_{pj}^2 \right)^2. \end{aligned}$$

The distance from  $\mathbf{a}_p$  to  $\mathcal{T}$  is  $O(\varrho \varepsilon^2)$  [15, Lemma 6], so  $\sum_{j=m+1}^d a_{pj}^2 = O(\varrho^2 \varepsilon^4)$ . It follows that  $\|\tilde{\mathbf{b}}_{p^*}\| \leq \sqrt{\varrho^2 \varepsilon^2 \cdot O(\varrho^2 \varepsilon^4) + O(\varrho^4 \varepsilon^8)} = O(\varrho^2 \varepsilon^3)$ . Since  $\|\mathbf{w}\| = 1$ ,

$$\mathbf{w}^t \cdot \tilde{\mathbf{b}}_{p^*} \leq \|\mathbf{w}\| \cdot \|\tilde{\mathbf{b}}_{p^*}\| = O(\varrho^2 \varepsilon^3).$$

The squared length of the projection of  $\mathbf{b}_{p^*}$  onto the support line of  $\mathbf{e}$  is

$$\left( \frac{\mathbf{e}^t \cdot \mathbf{b}_{p^*}}{\|\mathbf{e}\|} \right)^2 = \left( \frac{\mathbf{v}^t \cdot \bar{\mathbf{b}}_{p^*} + \mathbf{w}^t \cdot \tilde{\mathbf{b}}_{p^*}}{\sqrt{\|\mathbf{v}\|^2 + \|\mathbf{w}\|^2}} \right)^2 = \left( \frac{O(\varrho^2 \varepsilon^3)}{\sqrt{1 + O(\varepsilon^2)}} \right)^2 = O(\varrho^4 \varepsilon^6).$$

Summing this bound over the  $n$  rows of  $\mathbf{B}$  gives  $O(n \varrho^4 \varepsilon^6)$ , which is an upper bound of  $\sigma$ .  $\square$

## 7 Eigenvalues of $\mathbf{H}\mathbf{H}^t$

### 7.1 Preliminaries

Given a  $k \times l$  matrix  $\mathbf{U}$ , what happens to  $\|\mathbf{U}\|$  if we multiply each row  $\mathbf{u}_{i^*}$  by a factor which is less than or equal to one? Let  $\mathbf{V}$  be the matrix obtained after changing  $\mathbf{U}$ . Then

$$\begin{aligned} \|\mathbf{V}\| &= \left( \max_{\mathbf{e} \in \mathbb{R}^l, \|\mathbf{e}\|=1} \mathbf{e}^t \mathbf{V}^t \mathbf{V} \mathbf{e} \right)^{1/2} = \left( \max_{\mathbf{e} \in \mathbb{R}^l, \|\mathbf{e}\|=1} \sum_{i=1}^k (\mathbf{v}_{i^*} \mathbf{e})^2 \right)^{1/2} \\ &\leq \left( \max_{\mathbf{e} \in \mathbb{R}^l, \|\mathbf{e}\|=1} \sum_{i=1}^k (\mathbf{u}_{i^*} \mathbf{e})^2 \right)^{1/2} = \left( \max_{\mathbf{e} \in \mathbb{R}^l, \|\mathbf{e}\|=1} \mathbf{e}^t \mathbf{U}^t \mathbf{U} \mathbf{e} \right)^{1/2} \\ &= \|\mathbf{U}\|. \end{aligned}$$

Also, for every orthonormal  $k \times k$  matrix  $\mathbf{R}$ ,  $\mathbf{R}^t \mathbf{R} = \mathbf{I}_k$ , and therefore,

$$\|\mathbf{R}\mathbf{U}\| = \left( \max_{\mathbf{e} \in \mathbb{R}^k, \|\mathbf{e}\|=1} \mathbf{e}^t (\mathbf{R}\mathbf{U})^t (\mathbf{R}\mathbf{U}) \mathbf{e} \right)^{1/2} = \left( \max_{\mathbf{e} \in \mathbb{R}^k, \|\mathbf{e}\|=1} \mathbf{e}^t \mathbf{U}^t \mathbf{U} \mathbf{e} \right)^{1/2} = \|\mathbf{U}\|.$$

We need three more technical results concerning the angles between vectors and spaces. The first one is from [12]. The other two are folklore and we include their proofs for completeness.

**Lemma 7.1** ([12, Lemma 1.1]) *Let  $\mathbf{M}$  be an  $s \times s$  real symmetric matrix with eigenvalues  $\mu_1, \dots, \mu_s$  in an arbitrary order. Let  $\mathbf{v}_i$  denote a unit eigenvector of  $\mathbf{M}$  corresponding to  $\mu_i$ . If  $\mathbf{M} + \mathbf{M}'$  is a real symmetric matrix,  $\mu'$  is an eigenvalue of  $\mathbf{M} + \mathbf{M}'$ , and  $\mathbf{z}$  is a unit eigenvector of  $\mathbf{M} + \mathbf{M}'$  corresponding to  $\mu'$ , then for every  $r \in [1, s - 1]$ , the angle between  $\mathbf{z}$  and the space spanned by  $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$  is at most  $\arcsin\left(\frac{\|\mathbf{M}'\|}{\min_{i \in [r+1, s]} |\mu_i - \mu'|}\right)$ .*

**Lemma 7.2** *Let  $(\mathbf{U} \ \mathbf{V})$  be an  $s \times s$  orthonormal matrix such that  $\mathbf{U}$  is  $s \times (s - r)$  and  $\mathbf{V}$  is  $s \times r$  for some  $r \in [1, s - 1]$ . Let  $\mathbf{Z} = (\mathbf{z}_{*1} \ \dots \ \mathbf{z}_{*r})$  be an  $s \times r$  orthonormal matrix. The angle between the column spaces of  $\mathbf{V}$  and  $\mathbf{Z}$  is  $\arcsin(\|\mathbf{U}^t \mathbf{Z}\|) \leq \arcsin(\sqrt{r} \max_{i \in [1, r]} \|\mathbf{U}^t \mathbf{z}_{*i}\|)$ .*

*Proof.* Let  $\theta$  denote the angle between the column spaces of  $\mathbf{V}$  and  $\mathbf{Z}$ . For every unit vector  $\mathbf{z}$  in the column space of  $\mathbf{Z}$ , the sine of the angle between  $\mathbf{z}$  and the column space of  $\mathbf{V}$  is  $\|\mathbf{U}^t \mathbf{z}\|/\|\mathbf{z}\| = \|\mathbf{U}^t \mathbf{z}\|$ . Each such vector  $\mathbf{z}$  is a linear combination of the columns of  $\mathbf{Z}$ , i.e.,  $\mathbf{z} = \mathbf{Z}\mathbf{e}$  for some  $r \times 1$  unit vector  $\mathbf{e}$ . It follows that  $\sin \theta = \max_{\|\mathbf{e}\|=1} \|\mathbf{U}^t \mathbf{Z}\mathbf{e}\| = \|\mathbf{U}^t \mathbf{Z}\|$ . Moreover,  $\|\mathbf{U}^t \mathbf{Z}\| \leq \|\mathbf{U}^t \mathbf{Z}\|_F = \sqrt{\sum_{i=1}^r \|\mathbf{U}^t \mathbf{z}_{*i}\|^2} \leq \sqrt{r} \cdot \max_{i \in [1, r]} \|\mathbf{U}^t \mathbf{z}_{*i}\|$ .  $\square$

**Lemma 7.3** *Let  $\mathbf{v}$  be a  $r$ -dimensional vector. Let  $\{\mathbf{e}_i : 1 \leq i \leq k\}$  be an orthonormal basis of the column space of a  $r \times l$  matrix  $\mathbf{U}$ . If  $\theta_i$  denotes the acute angle between the support lines of  $\mathbf{v}$  and  $\mathbf{e}_i$ , then the acute angle between  $\mathbf{v}$  and the column space of  $\mathbf{U}$  is at least  $\arccos(\sum_{i=1}^k \cos \theta_i) = \arccos(\frac{1}{\|\mathbf{v}\|} \sum_{i=1}^k |\mathbf{v}^t \mathbf{e}_i|)$ , provided that  $\frac{1}{\|\mathbf{v}\|} \sum_{i=1}^k |\mathbf{v}^t \mathbf{e}_i| \leq 1$ .*

*Proof.* Take the projection of  $\mathbf{v}$  onto the column space of  $\mathbf{U}$ , and normalize the projection to a unit vector  $\mathbf{w}$ . So  $\mathbf{w} = \sum_{i=1}^k \alpha_i \mathbf{e}_i$  for some  $\alpha_i \in [-1, 1]$ . The cosine of the angle between  $\mathbf{v}$  and  $\mathbf{w}$  is equal to  $\frac{1}{\|\mathbf{v}\|} \mathbf{v}^t \mathbf{w} = \frac{1}{\|\mathbf{v}\|} \sum_{i=1}^k \alpha_i \mathbf{v}^t \mathbf{e}_i \leq \frac{1}{\|\mathbf{v}\|} \sum_{i=1}^k |\alpha_i| |\mathbf{v}^t \mathbf{e}_i| \leq \frac{1}{\|\mathbf{v}\|} \sum_{i=1}^k |\mathbf{v}^t \mathbf{e}_i|$ .  $\square$

The main result of this section is Lemma 7.10 which gives bounds on the eigenvalues of  $\mathbf{H}\mathbf{H}^t$ . The proof of Lemma 7.10 will be facilitated by a rotation of  $\mathbb{R}^d$  so that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . We prove below that such a rotation does not change the eigenvalues of  $\mathbf{H}\mathbf{H}^t$ .

**Lemma 7.4** *If we apply a rotation to  $\mathbb{R}^d$  that keeps the origin fixed, the eigenvalues of  $\mathbf{H}\mathbf{H}^t$  are preserved and the eigenvectors of  $\mathbf{H}\mathbf{H}^t$  are rotated correspondingly.*

*Proof.* Recall that  $\mathbf{L}\mathbf{A}\mathbf{R}^t$  denotes the thin SVD of  $\mathbf{B}$  and that  $\mathbf{H}\mathbf{H}^t = \mathbf{A}^t \mathbf{L} \mathbf{\Sigma} \mathbf{L}^t \mathbf{A}$  by Lemma 4.1(iii). Let  $\mathbf{M}$  be a  $d \times d$  rotation matrix. The proof of Lemma 6.2 reveals that the effect of applying  $\mathbf{M}$  is produced by an isometry in  $\mathbb{R}^{d_0}$ . It follows that the application of  $\mathbf{M}$  only changes the matrix  $\mathbf{R}$  but not  $\mathbf{L}$  or  $\mathbf{A}$  in the thin SVD of  $\mathbf{B}$ . Therefore, when we apply  $\mathbf{M}$ , the middle part  $\mathbf{L} \mathbf{\Sigma} \mathbf{L}^t$  of  $\mathbf{H}\mathbf{H}^t$  remains fixed and  $\mathbf{A}$  is changed to  $\mathbf{A}\mathbf{M}^t$  by the rotation. This changes  $\mathbf{H}\mathbf{H}^t$  to  $\mathbf{M}\mathbf{A}^t \mathbf{L} \mathbf{\Sigma} \mathbf{L}^t \mathbf{A}\mathbf{M}^t = \mathbf{M}\mathbf{H}\mathbf{H}^t \mathbf{M}^t$ . Since  $\mathbf{M}$  is a rotation matrix, multiplying  $\mathbf{M}$  on the left and  $\mathbf{M}^t$  on the right does not change the eigenvalues of  $\mathbf{H}\mathbf{H}^t$ , but it does rotate the eigenvectors of  $\mathbf{H}\mathbf{H}^t$  correspondingly.  $\square$

## 7.2 Analysis

Suppose that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . Then,  $\mathbf{A}$  can be divided into two submatrices. Let  $\mathbf{A}_T$  be the leftmost  $n \times m$  submatrix of  $\mathbf{A}$ . Let  $\mathbf{A}_N$  be the rightmost  $n \times (d - m)$  submatrix

of  $\mathbf{A}$ . That is,  $\mathbf{A} = (\mathbf{A}_T \ \mathbf{A}_N)$ . Lemma 4.1(iii) implies that  $\mathbf{H} = (\widehat{\mathbf{B}}^\dagger \mathbf{A})^t = \begin{pmatrix} \widehat{\mathbf{B}}^\dagger \mathbf{A}_T & \widehat{\mathbf{B}}^\dagger \mathbf{A}_N \end{pmatrix}^t$ . We define two submatrices of  $\mathbf{H}$  as follows.

$$\mathbf{H} = \begin{pmatrix} \mathbf{H}_T \\ \mathbf{H}_N \end{pmatrix}, \text{ where } \mathbf{H}_T \stackrel{\text{def}}{=} (\widehat{\mathbf{B}}^\dagger \mathbf{A}_T)^t \text{ and } \mathbf{H}_N \stackrel{\text{def}}{=} (\widehat{\mathbf{B}}^\dagger \mathbf{A}_N)^t.$$

Note that  $\mathbf{H}_T$  is an  $m \times d_0$  matrix and  $\mathbf{H}_N$  is a  $(d - m) \times d_0$  matrix. Our analysis begins with bounding  $\|\mathbf{H}_N\|$ .

**Lemma 7.5** *Assume that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ ,  $\|\mathbf{H}_N\| = O(\sqrt{n}\varrho\varepsilon^3/\lambda_{m_0+1})$ .*

*Proof.* Recall that  $f_\ell$ ,  $\ell \in [m + 1, d]$ , denotes a coordinate function of  $\mathcal{M}$  at the origin. Define:

$$\begin{aligned} \tilde{a}_{p\ell} &\stackrel{\text{def}}{=} \frac{1}{2} \mathbf{D}^2 f_\ell|_0 \left( (a_{p1} \ \cdots \ a_{pm})^t, (a_{p1} \ \cdots \ a_{pm})^t \right) \\ \tilde{\mathbf{A}}_N &\stackrel{\text{def}}{=} \begin{pmatrix} \tilde{a}_{1,m+1} & \cdots & \tilde{a}_{1d} \\ \vdots & \ddots & \vdots \\ \tilde{a}_{n,m+1} & \cdots & \tilde{a}_{nd} \end{pmatrix} \\ \tilde{\mathbf{H}}_N &\stackrel{\text{def}}{=} (\widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N)^t \end{aligned}$$

In the Taylor expansion of  $f_\ell$ , there is no constant term or first order term because  $\mathcal{M}$  contains the origin and the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . As a result, if  $\varepsilon$  is sufficiently small, then  $\tilde{a}_{p\ell}$  is close to  $a_{p\ell}$ , and therefore,  $\tilde{\mathbf{A}}_N$  and  $\tilde{\mathbf{H}}_N$  are approximations of  $\mathbf{A}_N$  and  $\mathbf{H}_N$ , respectively. Substituting  $\mathbf{A}_N = \tilde{\mathbf{A}}_N + (\mathbf{A}_N - \tilde{\mathbf{A}}_N)$  into  $\mathbf{H}_N = (\widehat{\mathbf{B}}^\dagger \mathbf{A}_N)^t$ , we obtain  $\mathbf{H}_N = \tilde{\mathbf{H}}_N + (\widehat{\mathbf{B}}^\dagger (\mathbf{A}_N - \tilde{\mathbf{A}}_N))^t$ . Therefore,  $\|\mathbf{H}_N\| \leq \|\tilde{\mathbf{H}}_N\| + \|\widehat{\mathbf{B}}^\dagger (\mathbf{A}_N - \tilde{\mathbf{A}}_N)\|$ .

We first bound  $\|\widehat{\mathbf{B}}^\dagger (\mathbf{A}_N - \tilde{\mathbf{A}}_N)\|$ . For every  $p \in [1, n]$ , the  $p$ th row of  $\mathbf{A}_N - \tilde{\mathbf{A}}_N$  contains the third and higher order terms in the Taylor expansions of  $f_\ell((a_{p1} \ \cdots \ a_{pm})^t)$  for  $\ell \in [m + 1, d]$ . By Lemma 5.1(ii), if  $\varepsilon$  is small enough, the 2-norm of each row of  $\mathbf{A}_N - \tilde{\mathbf{A}}_N$  is  $O(\varrho\varepsilon^3)$ . Thus,  $\|\mathbf{A}_N - \tilde{\mathbf{A}}_N\| \leq \|\mathbf{A}_N - \tilde{\mathbf{A}}_N\|_F = O(\sqrt{n}\varrho\varepsilon^3)$ . Since  $\lambda_{m_0+1}$  is the smallest singular value of  $\widehat{\mathbf{B}}$ ,  $\|\widehat{\mathbf{B}}^\dagger\| = 1/\lambda_{m_0+1}$ . Therefore,  $\|\widehat{\mathbf{B}}^\dagger (\mathbf{A}_N - \tilde{\mathbf{A}}_N)\| \leq \|\widehat{\mathbf{B}}^\dagger\| \cdot \|\mathbf{A}_N - \tilde{\mathbf{A}}_N\| = O(\sqrt{n}\varrho\varepsilon^3/\lambda_{m_0+1})$ .

It remains to show that  $\|\tilde{\mathbf{H}}_N\| = O(\sqrt{n}\varrho\varepsilon^3/\lambda_{m_0+1})$ . Observe that  $\|\tilde{\mathbf{H}}_N\| = \|\tilde{\mathbf{H}}_N^t\| = \|\widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N\|$ . The smallest singular value of  $\mathbf{B}$  is smaller than or equal to the smallest singular value of  $\widehat{\mathbf{B}}$  by construction. Intuitively, we would expect  $\|\widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N\| \leq \|\mathbf{B}^\dagger \tilde{\mathbf{A}}_N\|$ . Therefore, we can bound  $\|\widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N\|$  if we can bound  $\|\mathbf{B}^\dagger \tilde{\mathbf{A}}_N\|$ .

Recall that  $\mathbf{D}^2 f_\ell|_0$  can be viewed as the  $m \times m$  symmetric matrix  $(\partial^2 f_\ell / \partial x_i \partial x_j)_{i,j \in [1,m]}$  with every entry evaluated at the origin. Let  $q_{\ell,ij}$  denote the  $(i, j)$  entry of  $\mathbf{D}^2 f_\ell|_0$ . By definition,  $\tilde{a}_{p\ell} = \frac{1}{2} (a_{p1} \ \cdots \ a_{pm}) \cdot \mathbf{D}^2 f_\ell|_0 \cdot (a_{p1} \ \cdots \ a_{pm})^t$ . Expanding this equation gives:

$$\tilde{a}_{p\ell} = \sum_{i=1}^m \sum_{j=1}^m \frac{1}{2} q_{\ell,ij} a_{pi} a_{pj} = \sum_{i=1}^m \frac{1}{2} q_{\ell,ii} a_{pi}^2 + \sum_{i=1}^m \sum_{j=i+1}^m q_{\ell,ij} a_{pi} a_{pj}.$$

We extend the range of  $i$  and  $j$  in  $q_{\ell,ij}$  to  $[1, d]$  by letting  $q_{\ell,ij} = 0$  whenever  $i \in [m + 1, d]$  or  $j \in [m + 1, d]$ . Then, define  $\mathbf{Z}$  to be the  $d_0 \times (d - m)$  matrix  $(\mathbf{z}_{*1} \ \cdots \ \mathbf{z}_{*(d-m)})$ , where

$$\mathbf{z}_{*\ell-m} = \left( \frac{1}{\sqrt{2}} q_{\ell,11} \quad q_{\ell,12} \quad \cdots \quad q_{\ell,1d} \quad \frac{1}{\sqrt{2}} q_{\ell,22} \quad q_{\ell,23} \quad \cdots \quad q_{\ell,2d} \quad \cdots \quad \frac{1}{\sqrt{2}} q_{\ell,dd} \right)^t.$$

The definition of  $\mathbf{z}_{*\ell-m}$  is crafted so that  $\mathbf{B} \mathbf{z}_{*\ell-m}$  equals  $(\tilde{a}_{1\ell} \ \cdots \ \tilde{a}_{n\ell})^t$ . Therefore,

$$\mathbf{B} \mathbf{Z} = \tilde{\mathbf{A}}_N. \tag{7.1}$$

Consider the square matrix obtained by appending  $d_0 - d + m$  zero columns to the right of  $\mathbf{Z}$ . Then,  $\|\mathbf{Z}\|$  is equal to the 2-norm of this square matrix, which by the Gershgorin Circle Theorem is at most  $\max_{j \in [1, d-m]} \sum_{i=1}^{d_0} |z_{ij}|$ . Among the entries in  $\mathbf{z}_{*j}$ ,  $m_0$  entries are from the upper triangular portion of  $\mathbf{D}^2 f_{j+m}|_0$  and the other  $d_0 - m_0$  entries are zeros by definition. Thus,  $\sum_{i=1}^{d_0} |z_{ij}| \leq \sqrt{m_0} \cdot \|\mathbf{z}_{*j}\|$ . By Lemma 5.1(i),  $\|\mathbf{D}^2 f_{j+m}|_0\| = O(1/\varrho)$  and therefore,  $\|\mathbf{z}_{*j}\| \leq \|\mathbf{D}^2 f_{j+m}|_0\|_F \leq \sqrt{m} \cdot \|\mathbf{D}^2 f_{j+m}|_0\| = O(1/\varrho)$ . By the Gershgorin Circle Theorem,  $\|\mathbf{Z}\| \leq \max_{j \in [1, d-m]} \sum_{i=1}^{d_0} |z_{ij}| \leq \max_{j \in [1, d-m]} \sqrt{m_0} \cdot \|\mathbf{z}_{*j}\| = O(1/\varrho)$ . Since  $d_0 - m_0$  rows of  $\mathbf{Z}$  contain only zeros by construction,  $\|\mathbf{Z}\|_F \leq \sqrt{m_0} \cdot \|\mathbf{Z}\| = O(1/\varrho)$ .

$\mathbf{B}^\dagger \tilde{\mathbf{A}}_N$  also satisfies (7.1), i.e.,  $\mathbf{B}(\mathbf{B}^\dagger \tilde{\mathbf{A}}_N) = \tilde{\mathbf{A}}_N$ . By the property of pseudoinverse,  $\|\mathbf{B}^\dagger \tilde{\mathbf{A}}_N\|_F$  is no more than  $\|\mathbf{Z}\|_F$  or the Frobenius norm of any matrix that satisfies (7.1). As a result,

$$\|\mathbf{B}^\dagger \tilde{\mathbf{A}}_N\| \leq \|\mathbf{B}^\dagger \tilde{\mathbf{A}}_N\|_F \leq \|\mathbf{Z}\|_F = O(1/\varrho). \quad (7.2)$$

We relate  $\|\widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N\|$  to  $\|\mathbf{B}^\dagger \tilde{\mathbf{A}}_N\|$  as follows. Recall that  $\mathbf{L}\mathbf{\Lambda}\mathbf{R}^t$  and  $\widehat{\mathbf{L}}\widehat{\mathbf{\Lambda}}\widehat{\mathbf{R}}^t$  are the thin SVDs of  $\mathbf{B}$  and  $\widehat{\mathbf{B}}$ , respectively, and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$  are the singular values of  $\mathbf{B}$ . Let  $\sigma_i$  and  $\widehat{\sigma}_i$  denote the diagonal entries of  $\mathbf{\Lambda}^\dagger$  and  $\widehat{\mathbf{\Lambda}}^\dagger$ , respectively. By Lemmas 6.5 and 6.6, it holds with probability  $1 - O(n^{-1/3})$  that  $\lambda_{m_0} = \Theta(\sqrt{n}\varrho^2\varepsilon^2)$  and  $\lambda_{m_0+1} = O(\sqrt{n}\varrho^2\varepsilon^3)$ . Therefore,

$$\begin{aligned} \forall i \in [1, m_0], \quad \widehat{\sigma}_i &= \sigma_i = 1/\lambda_i, \\ \forall i \in [m_0 + 1, n], \quad \widehat{\sigma}_i &= \frac{1}{\lambda_{m_0+1}} \leq \frac{1}{\lambda_i} = \sigma_i. \end{aligned}$$

(For simplicity, we assume that the smallest singular value  $\lambda_n$  of  $\mathbf{\Lambda}$  is positive, and therefore,  $\widehat{\mathbf{\Lambda}}$  has  $n$  positive diagonal entries. Otherwise, if  $\lambda_i = 0$ , both  $\sigma_i$  and  $\widehat{\sigma}_i$  are zero by our definition of  $\widehat{\mathbf{\Lambda}}$ .)

The full SVD of  $\mathbf{B}$  is  $\mathbf{L}(\mathbf{\Lambda} \ \mathbf{0}_{n, d_0-n})\bar{\mathbf{R}}^t$ , where  $\bar{\mathbf{R}}$  consists of the  $d_0$  unit eigenvectors of  $\mathbf{B}^t\mathbf{B}$  and  $\mathbf{R}$  is the leftmost  $d_0 \times n$  submatrix of  $\bar{\mathbf{R}}$ . It follows that  $\mathbf{L}(\widehat{\mathbf{\Lambda}} \ \mathbf{0}_{n, d_0-n})\bar{\mathbf{R}}^t$  is the full SVD of  $\widehat{\mathbf{B}}$ , and therefore,  $\mathbf{B}^\dagger = \bar{\mathbf{R}}(\mathbf{\Lambda}^\dagger \ \mathbf{0}_{n, d_0-n})^t \mathbf{L}^t$  and  $\widehat{\mathbf{B}}^\dagger = \bar{\mathbf{R}}(\widehat{\mathbf{\Lambda}}^\dagger \ \mathbf{0}_{n, d_0-n})^t \mathbf{L}^t$ . Observe that:

- For  $i \in [1, n]$ , the  $i$ th row in  $\bar{\mathbf{R}}^t \widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N$  equals the  $i$ th row in  $\bar{\mathbf{R}}^t \mathbf{B}^\dagger \tilde{\mathbf{A}}_N$  multiplied by  $\widehat{\sigma}_i/\sigma_i$ .
- The bottom  $d_0 - n$  rows in both  $\bar{\mathbf{R}}^t \widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N$  and  $\bar{\mathbf{R}}^t \mathbf{B}^\dagger \tilde{\mathbf{A}}_N$  contain only zeros.

Therefore,  $\|\bar{\mathbf{R}}^t \widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N\| \leq \|\bar{\mathbf{R}}^t \mathbf{B}^\dagger \tilde{\mathbf{A}}_N\|$  as  $\widehat{\sigma}_i \leq \sigma_i$  for  $i \in [1, n]$ .

Multiplying  $\tilde{\mathbf{H}}_N^t$ ,  $\widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N$  and  $\mathbf{B}^\dagger \tilde{\mathbf{A}}_N$  on the left by  $\bar{\mathbf{R}}^t$  does not change their 2-norms. Therefore,

$$\|\tilde{\mathbf{H}}_N\| = \|\tilde{\mathbf{H}}_N^t\| = \|\bar{\mathbf{R}}^t \widehat{\mathbf{B}}^\dagger \tilde{\mathbf{A}}_N\| \leq \|\bar{\mathbf{R}}^t \mathbf{B}^\dagger \tilde{\mathbf{A}}_N\| = \|\mathbf{B}^\dagger \tilde{\mathbf{A}}_N\| \stackrel{(7.2)}{=} O(1/\varrho).$$

Thus,  $\|\tilde{\mathbf{H}}_N\| = O(\sqrt{n}\varrho\varepsilon^3/\lambda_{m_0+1})$  because  $\lambda_{m_0+1} = O(\sqrt{n}\varrho^2\varepsilon^3)$  by Lemma 6.6.  $\square$

Before we analyze the singular values of  $\mathbf{H}_T$ , we prove two technical results about the column vectors of  $\mathbf{A}_T$  and the column space of  $\mathbf{B}_{TT}$ . First, we show that the acute angle between the support lines of any two distinct columns of  $\mathbf{A}_T$  is large and  $\mathbf{A}_T$  has rank  $m$ .

**Lemma 7.6** *Assume that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , for every distinct  $i, j \in [1, m]$ , the acute angle between the support lines of  $\mathbf{a}_{*i}$  and  $\mathbf{a}_{*j}$  is at least  $\arccos(\frac{1}{9}m^{-1}n^{-1/3})$ , and  $\mathbf{A}_T$  has rank  $m$ .*

*Proof.* Take an arbitrary pair of distinct columns  $\mathbf{a}_{*i}$  and  $\mathbf{a}_{*j}$  of  $\mathbf{A}_T$ . By Lemma 5.4(iii), it holds with probability  $1 - O(n^{-1/3})$  that  $\|\mathbf{a}_{*i}\|$  and  $\|\mathbf{a}_{*j}\|$  are at least  $c_1\sqrt{n}\varrho\varepsilon$  for some constant  $c_1 > 0$ . The inner product  $\mathbf{a}_{*i}^t \cdot \mathbf{a}_{*j}$  equals  $\sum_{p=1}^n a_{pi}a_{pj}$ . Apply Lemma 5.4(vi) with the

constant  $c = \frac{1}{18}c_1^2m^{-1}$ . It implies that if  $\varepsilon$  is small enough, then with probability  $1 - O(n^{-1/3})$ ,  $|\mathbf{a}_{*i}^t \cdot \mathbf{a}_{*j}| \leq cn^{2/3}\varrho^2\varepsilon^2 + O(n\varrho^2\varepsilon^4) \leq 2cn^{2/3}\varrho^2\varepsilon^2 = \frac{c_1^2}{9}m^{-1}n^{2/3}\varrho^2\varepsilon^2$ . Thus,  $|\mathbf{a}_{*i}^t \cdot \mathbf{a}_{*j}| / (\|\mathbf{a}_{*i}\| \|\mathbf{a}_{*j}\|) \leq \frac{1}{9}m^{-1}n^{-1/3}$ , which implies that the acute angle between the support lines of  $\mathbf{a}_{*i}$  and  $\mathbf{a}_{*j}$  is at least  $\arccos(\frac{1}{9}m^{-1}n^{-1/3})$ .

We prove that  $\mathbf{A}_T$  has rank  $m$  by bounding the minimum eigenvalue of  $\mathbf{A}_T^t \mathbf{A}_T$  away from zero. Observe that:

$$\mathbf{A}_T^t \mathbf{A}_T = \begin{pmatrix} \sum_{p=1}^n a_{p1}^2 & \sum_{p=1}^n a_{p1}a_{p2} & \cdots & \sum_{p=1}^n a_{p1}a_{pm} \\ \sum_{p=1}^n a_{p2}a_{p1} & \sum_{p=1}^n a_{p2}^2 & \cdots & \sum_{p=1}^n a_{p2}a_{pm} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{p=1}^n a_{pm}a_{p1} & \sum_{p=1}^n a_{pm}a_{p2} & \cdots & \sum_{p=1}^n a_{pm}^2 \end{pmatrix}$$

By the Gershgorin Circle Theorem, the minimum eigenvalue of  $\mathbf{A}_T^t \mathbf{A}_T$  is greater than or equal to  $\min_{i \in [1, m]} \left\{ \sum_{p=1}^n a_{pi}^2 - \sum_{j \in [1, m] \setminus \{i\}} \left| \sum_{p=1}^n a_{pi}a_{pj} \right| \right\}$ . Apply Lemma 5.4(iii) with  $c = \frac{1}{3(m+2)}$  to obtain  $\sum_{p=1}^n a_{pi}^2 \geq \frac{1}{m+2}n\varrho^2\varepsilon^2 - \frac{1}{3(m+2)}n^{2/3}\varrho^2\varepsilon^2 - O(n\varrho^2\varepsilon^4)$  with probability  $1 - O(n^{-1/3})$ . Apply Lemma 5.4(vi) with  $c = \frac{1}{3(m-1)(m+2)}$  to obtain  $\sum_{j \in [1, m] \setminus \{i\}} \left| \sum_{p=1}^n a_{pi}a_{pj} \right| \leq \frac{1}{3(m+2)}n^{2/3}\varrho^2\varepsilon^2 + O(n\varrho^2\varepsilon^4)$  with probability  $1 - O(n^{-1/3})$ . Then, with probability  $1 - O(n^{-1/3})$ , the minimum eigenvalue of  $\mathbf{A}_T^t \mathbf{A}_T$  is at least  $\frac{1}{3(m+2)}n\varrho^2\varepsilon^2 - O(n\varrho^2\varepsilon^4)$ , which is positive when  $\varepsilon$  is small enough.  $\square$

Next, we show a lower bound on the angle between any column vector of  $\mathbf{A}_T$  and the column space of  $\mathbf{B}_{TT}$ .

**Lemma 7.7** *Assume that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , every column vector of  $\mathbf{A}_T$  makes an angle  $\arccos(\frac{1}{9}m^{-1}n^{-1/3})$  or more with the column space of  $\mathbf{B}_{TT}$ .*

*Proof.* We first introduce three constants  $c_1, c_2$  and  $c_3$ . It holds with probability  $1 - O(n^{-1/3})$  that all eigenvalues of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  are greater than  $c_1 n \varrho^4 \varepsilon^4$  for some constant  $c_1 > 0$  (Lemma 6.3), the 2-norm of each column of  $\mathbf{B}_{TT}$  is greater than  $c_2 \sqrt{n} \varrho^2 \varepsilon^2$  for some constant  $c_2 > 0$  (Lemma 5.4(i) and (ii)), and the 2-norm of every column of  $\mathbf{A}_T$  is greater than  $c_3 \sqrt{n} \varrho \varepsilon$  for some constant  $c_3 > 0$  (Lemma 5.4(iii)).

Let  $\mathbf{e}_1, \dots, \mathbf{e}_{m_0}$  be  $m_0$  unit eigenvectors corresponding to the eigenvalues of  $\mathbf{B}_{TT} \mathbf{B}_{TT}^t$ . Let  $\mathbf{b}_{*1}, \mathbf{b}_{*2}, \dots, \mathbf{b}_{*m_0}$  denote the columns of  $\mathbf{B}_{TT}$ .

We first show that for every  $j \in [1, m_0]$ ,  $\mathbf{e}_j = \sum_{i=1}^{m_0} \beta_{ij} \mathbf{b}_{*i} / \|\mathbf{b}_{*i}\|$  for some coefficients  $\beta_{ij}$ 's such that  $|\beta_{ij}| \leq 1/\sqrt{c_1}$ . By Lemma 6.3,  $\mathbf{B}_{TT} \mathbf{B}_{TT}^t$  has rank  $m_0$ , implying that  $\mathbf{e}_j = \sum_{i=1}^{m_0} \beta_{ij} \mathbf{b}_{*i} / \|\mathbf{b}_{*i}\|$  for some coefficients  $\beta_{ij}$ 's. It remains to bound  $|\beta_{ij}|$ . We put the equations  $\mathbf{e}_j = \sum_{i=1}^{m_0} \beta_{ij} \mathbf{b}_{*i} / \|\mathbf{b}_{*i}\|$  in matrix form as follows.

$$\mathbf{E} = (\mathbf{e}_1 \quad \cdots \quad \mathbf{e}_{m_0}) = \mathbf{B}_{TT} \begin{pmatrix} \frac{1}{\|\mathbf{b}_{*1}\|} \beta_{11} & \cdots & \frac{1}{\|\mathbf{b}_{*1}\|} \beta_{1m_0} \\ \vdots & \ddots & \vdots \\ \frac{1}{\|\mathbf{b}_{*m_0}\|} \beta_{m_01} & \cdots & \frac{1}{\|\mathbf{b}_{*m_0}\|} \beta_{m_0m_0} \end{pmatrix} \quad (7.3)$$

The thin SVD of  $\mathbf{B}_{TT}^t$  is  $\mathbf{V} \mathbf{D}^t$ , where  $\mathbf{D}$  is an  $m_0 \times m_0$  diagonal matrix whose  $(j, j)$  entry is the square root  $\mu_j$  of the  $j$ -th largest eigenvalue of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$ , and  $\mathbf{V}$  is an  $m_0 \times m_0$  matrix whose  $j$ -th column  $\mathbf{v}_{*j}$  is the unit eigenvector of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  corresponding to  $\mu_j^2$ . Multiplying both sides

of (7.3) by  $\text{VD}^{-1}\mathbf{E}^t$  gives:

$$\text{VD}^{-1} = \begin{pmatrix} \frac{1}{\|\mathbf{b}_{*1}\|}\beta_{11} & \cdots & \frac{1}{\|\mathbf{b}_{*1}\|}\beta_{1m_0} \\ \vdots & \ddots & \vdots \\ \frac{1}{\|\mathbf{b}_{*m_0}\|}\beta_{m_01} & \cdots & \frac{1}{\|\mathbf{b}_{*m_0}\|}\beta_{m_0m_0} \end{pmatrix}.$$

Comparing the matrices on the two sides term by term shows that  $\beta_{ij} = \|\mathbf{b}_{*i}\|v_{ij}/\mu_j$  for  $i, j \in [1, m_0]$ . By assumption,  $\mu_j \geq \sqrt{c_1}\sqrt{n}\varrho^2\varepsilon^2$ . Since  $|a_{pr}a_{ps}| \leq \varrho^2\varepsilon^2$  for  $r, s \in [1, d]$ ,  $\|\mathbf{b}_{*i}\| \leq \sqrt{n}\varrho^2\varepsilon^2$ . Also,  $|v_{ij}| \leq 1$  as  $\mathbf{v}_{*j}$  is a unit vector. It follows that  $|\beta_{ij}| \leq 1/\sqrt{c_1}$ .

We bound  $\mathbf{a}_{*k}^t \cdot \mathbf{b}_{*i}$  for  $k \in [1, m]$  and  $i \in [1, m_0]$  as follows. The inner product  $\mathbf{a}_{*k}^t \cdot \mathbf{b}_{*i}$  is equal to  $\sum_{p=1}^n \frac{1}{\sqrt{2}}a_{pk}a_{pr}^2$  or  $\sum_{p=1}^n a_{pk}a_{pr}a_{ps}$  for some possibly non-distinct  $k, r, s \in [1, m]$ . Apply Lemma 5.4(v) with the constant  $c = \frac{\sqrt{c_1}c_2c_3}{18}m^{-1}m_0^{-2}$ . It implies that, with probability  $1 - O(n^{-1/3})$ ,  $|\mathbf{a}_{*k}^t \cdot \mathbf{b}_{*i}| \leq cn^{2/3}\varrho^3\varepsilon^3 + O(n\varrho^3\varepsilon^5)$ , which is at most  $2cn^{2/3}\varrho^3\varepsilon^3 = \frac{\sqrt{c_1}c_2c_3}{9}m^{-1}m_0^{-2}n^{2/3}\varrho^3\varepsilon^3$  when  $\varepsilon$  is sufficiently small.

For  $k \in [1, m]$  and  $j \in [1, m_0]$ ,  $\frac{1}{\|\mathbf{a}_{*k}\|}|\mathbf{a}_{*k}^t \cdot \mathbf{e}_j| \leq \sum_{i=1}^{m_0} \frac{|\beta_{ij}|}{\|\mathbf{a}_{*k}\|\|\mathbf{b}_{*i}\|}|\mathbf{a}_{*k}^t \cdot \mathbf{b}_{*i}| \leq \sum_{i=1}^{m_0} \frac{1}{\sqrt{c_1}c_2c_3n\varrho^3\varepsilon^3} \cdot \frac{\sqrt{c_1}c_2c_3}{9}m^{-1}m_0^{-2}n^{2/3}\varrho^3\varepsilon^3 = \frac{1}{9}m^{-1}m_0^{-1}n^{-1/3}$ . Therefore,  $\frac{1}{\|\mathbf{a}_{*k}\|} \sum_{j=1}^{m_0} |\mathbf{a}_{*k}^t \cdot \mathbf{e}_j| \leq \frac{1}{9}m^{-1}n^{-1/3}$ . Then, Lemma 7.3 implies that the acute angle between  $\mathbf{a}_{*k}$  and the column space of  $\mathbf{B}_{TT}$  is  $\arccos(\frac{1}{9}m^{-1}n^{-1/3})$  or more.  $\square$

The matrix  $\widehat{\Lambda}$  in the thin SVD  $\widehat{\Lambda}\widehat{\mathbf{L}}\widehat{\mathbf{R}}^t$  of  $\widehat{\mathbf{B}}$  can be partitioned into blocks to separate the  $m_0$  largest singular values of  $\mathbf{B}$  from the  $\lambda_{m_0+1}$ 's. The matrices  $\mathbf{L}$  and  $\mathbf{R}$  can then be partitioned correspondingly. Specifically, we obtain

$$\begin{aligned} \widehat{\Lambda} &= \begin{pmatrix} \widehat{\Lambda}_0 & \mathbf{0}_{m_0, n-m_0} \\ \mathbf{0}_{n-m_0, m_0} & \lambda_{m_0+1}\mathbf{I}_{n-m_0} \end{pmatrix}, \\ \mathbf{L} &= \begin{pmatrix} \underbrace{\mathbf{L}_0}_{m_0 \text{ columns}} & \underbrace{\mathbf{L}_1}_{n-m_0 \text{ columns}} \end{pmatrix}, \\ \mathbf{R} &= \begin{pmatrix} \underbrace{\mathbf{R}_0}_{m_0 \text{ columns}} & \underbrace{\mathbf{R}_1}_{n-m_0 \text{ columns}} \end{pmatrix}. \end{aligned} \tag{7.4}$$

Note that  $\|\mathbf{L}_0\| = \|\mathbf{L}_1\| = 1$ .

Our subsequent analysis of the eigenvalues of  $\mathbf{H}_T$  requires an upper bound on the angle between the column spaces of  $\mathbf{L}_0$  and  $\mathbf{B}_{TT}$ , which is given in the following result.

**Lemma 7.8** *Assume that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , the angle between the column spaces of  $\mathbf{L}_0$  and  $\mathbf{B}_{TT}$  is  $O(\varepsilon^2)$ .*

*Proof.* Recall that  $\mathbf{B} = (\mathbf{B}_{TT} \ \mathbf{B}_{TN} \ \mathbf{B}_{NN})$ . Therefore,  $\mathbf{B}\mathbf{B}^t = \mathbf{B}_{TT}\mathbf{B}_{TT}^t + \mathbf{B}_{TN}\mathbf{B}_{TN}^t + \mathbf{B}_{NN}\mathbf{B}_{NN}^t$ . Let  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{m_0}$  be the  $m_0$  largest eigenvalues of  $\mathbf{B}_{TT}\mathbf{B}_{TT}^t$ . Let  $\mathbf{V}$  be an  $n \times m_0$  matrix whose columns are the unit eigenvectors of  $\mathbf{B}_{TT}\mathbf{B}_{TT}^t$  corresponding to  $\mu_1, \dots, \mu_{m_0}$ . The diagonalization of  $\mathbf{B}_{TT}\mathbf{B}_{TT}^t$  is:

$$\mathbf{B}_{TT}\mathbf{B}_{TT}^t = (\mathbf{V} \ *) \cdot \text{diag}_n(\mu_1, \dots, \mu_{m_0}, 0, \dots, 0) \cdot (\mathbf{V} \ *)^t. \tag{7.5}$$

Take any column vector  $\mathbf{z}$  of  $\mathbf{L}_0$ . Thus,  $\mathbf{B}\mathbf{B}^t\mathbf{z} = \mu'\mathbf{z}$  where  $\mu'$  is one of the largest  $m_0$  eigenvalues of  $\mathbf{B}\mathbf{B}^t$ . We expand  $\mathbf{B}\mathbf{B}^t$  to obtain  $(\mathbf{B}_{TT}\mathbf{B}_{TT}^t + \mathbf{B}_{TN}\mathbf{B}_{TN}^t + \mathbf{B}_{NN}\mathbf{B}_{NN}^t)\mathbf{z} = \mu'\mathbf{z}$ . We invoke Lemma 7.1 with  $\mathbf{M} = \mathbf{B}_{TT}\mathbf{B}_{TT}^t$  and  $\mathbf{M}' = \mathbf{B}_{TN}\mathbf{B}_{TN}^t + \mathbf{B}_{NN}\mathbf{B}_{NN}^t$ . Lemma 7.1 and (7.5) imply that the angle between  $\mathbf{z}$  and the column space of  $\mathbf{V}$  is at most  $\arcsin(\|\mathbf{M}'\|/|\mu'|)$ . By Lemma 6.4,

$\|M'\| \leq \|B_{TN}\|^2 + \|B_{NN}\|^2 = O(n\rho^4\varepsilon^6)$ . Lemma 6.5 implies that the largest  $m_0$  eigenvalues of  $BB^t$ , including  $\mu'$ , are  $\Theta(n\rho^4\varepsilon^4)$  with probability  $1 - O(n^{-1/3})$ . It follows that, with probability  $1 - O(n^{-1/3})$ , for every column  $\mathbf{z}$  of  $\mathbf{L}_0$ , the angle between  $\mathbf{z}$  and the column space of  $\mathbf{V}$  is at most  $\arcsin(O(\varepsilon^2)) = O(\varepsilon^2)$ .

Let  $\mathbf{U}$  be an  $n \times (n - m_0)$  matrix such that the columns of  $\mathbf{U}$  and  $\mathbf{V}$  form an orthonormal basis of the column space of  $\mathbf{B}$ . For every column  $\mathbf{z}$  of  $\mathbf{L}_0$ , since  $\mathbf{z}$  makes an  $O(\varepsilon^2)$  angle with the column space of  $\mathbf{V}$ , the angle between  $\mathbf{z}$  and the column space of  $\mathbf{U}$  is  $\pi/2 - O(\varepsilon^2)$ . That is,  $\|\mathbf{U}^t\mathbf{z}\| = \cos(\pi/2 - O(\varepsilon^2)) = O(\varepsilon^2)$ . Then, Lemma 7.2 implies that the angle between the column spaces of  $\mathbf{L}_0$  and  $\mathbf{V}$  is  $\arcsin(\sqrt{m_0} \cdot O(\varepsilon^2)) = O(\varepsilon^2)$ . Since the columns in  $\mathbf{V}$  form an orthonormal basis of the column space of  $\mathbf{B}_{TT}$ , the angle between the column spaces of  $\mathbf{L}_0$  and  $\mathbf{B}_{TT}$  is  $O(\varepsilon^2)$ .  $\square$

We are ready to bound the eigenvalues of  $\mathbf{H}_T$ . The analysis uses the tools that we have just developed, namely Lemmas 7.6, 7.7, and 7.8.

**Lemma 7.9** *Assume that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , the singular values of  $\mathbf{H}_T$  are  $\Theta(\sqrt{n}\rho\varepsilon/\lambda_{m_0+1})$  and so  $\|\mathbf{H}_T\| = \Theta(\sqrt{n}\rho\varepsilon/\lambda_{m_0+1})$ .*

*Proof.* Refer to the partitions of  $\widehat{\Lambda}$ ,  $\mathbf{L}$  and  $\mathbf{R}$  in (7.4). Note that  $\widehat{\Lambda}_0 = \text{diag}_{m_0}(\lambda_1, \dots, \lambda_{m_0})$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{m_0}$  are the  $m_0$  largest singular values of  $\mathbf{B}$ . Then,

$$\begin{aligned} \mathbf{H}_T^t &= \widehat{\mathbf{B}}^\dagger \mathbf{A}_T = (\mathbf{R}_0 \ \mathbf{R}_1) \begin{pmatrix} \widehat{\Lambda}_0^\dagger & \mathbf{0}_{m_0, n-m_0} \\ \mathbf{0}_{n-m_0, m_0} & \frac{1}{\lambda_{m_0+1}} \mathbf{I}_{n-m_0} \end{pmatrix} \begin{pmatrix} \mathbf{L}_0^t \\ \mathbf{L}_1^t \end{pmatrix} \mathbf{A}_T \\ &= \mathbf{R}_0 \widehat{\Lambda}_0^\dagger \mathbf{L}_0^t \mathbf{A}_T + \frac{1}{\lambda_{m_0+1}} \mathbf{R}_1 \mathbf{L}_1^t \mathbf{A}_T. \end{aligned}$$

Since every column vector of  $\mathbf{R}_0$  is orthogonal to any column vector of  $\mathbf{R}_1$  (i.e.,  $\mathbf{R}_0^t \mathbf{R}_1 = \mathbf{0}_{m_0, n-m_0}$  and  $\mathbf{R}_1^t \mathbf{R}_0 = \mathbf{0}_{n-m_0, m_0}$ ), we obtain

$$\begin{aligned} \mathbf{H}_T \mathbf{H}_T^t &= \left( \mathbf{R}_0 \widehat{\Lambda}_0^\dagger \mathbf{L}_0^t \mathbf{A}_T + \frac{1}{\lambda_{m_0+1}} \mathbf{R}_1 \mathbf{L}_1^t \mathbf{A}_T \right)^t \left( \mathbf{R}_0 \widehat{\Lambda}_0^\dagger \mathbf{L}_0^t \mathbf{A}_T + \frac{1}{\lambda_{m_0+1}} \mathbf{R}_1 \mathbf{L}_1^t \mathbf{A}_T \right) \\ &= \left( \mathbf{R}_0 \widehat{\Lambda}_0^\dagger \mathbf{L}_0^t \mathbf{A}_T \right)^t \left( \mathbf{R}_0 \widehat{\Lambda}_0^\dagger \mathbf{L}_0^t \mathbf{A}_T \right) + \frac{1}{\lambda_{m_0+1}^2} \left( \mathbf{R}_1 \mathbf{L}_1^t \mathbf{A}_T \right)^t \left( \mathbf{R}_1 \mathbf{L}_1^t \mathbf{A}_T \right) \\ &= \mathbf{A}_T^t \mathbf{L}_0 \left( \widehat{\Lambda}_0^\dagger \right)^2 \mathbf{L}_0^t \mathbf{A}_T + \frac{1}{\lambda_{m_0+1}^2} \mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T. \end{aligned}$$

The three matrices  $\mathbf{H}_T \mathbf{H}_T^t$ ,  $\mathbf{A}_T^t \mathbf{L}_0 \left( \widehat{\Lambda}_0^\dagger \right)^2 \mathbf{L}_0^t \mathbf{A}_T$ , and  $\frac{1}{\lambda_{m_0+1}^2} \mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T$  are symmetric and positive semi-definite. It follows that the maximum eigenvalue of  $\mathbf{H}_T \mathbf{H}_T^t$  is at most the sum of the maximum eigenvalues of  $\mathbf{A}_T^t \mathbf{L}_0 \left( \widehat{\Lambda}_0^\dagger \right)^2 \mathbf{L}_0^t \mathbf{A}_T$  and  $\frac{1}{\lambda_{m_0+1}^2} \mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T$ , and the minimum eigenvalue of  $\mathbf{H}_T \mathbf{H}_T^t$  is at least the maximum of the minimum eigenvalues of  $\mathbf{A}_T^t \mathbf{L}_0 \left( \widehat{\Lambda}_0^\dagger \right)^2 \mathbf{L}_0^t \mathbf{A}_T$  and  $\frac{1}{\lambda_{m_0+1}^2} \mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T$ .

For every  $i \in [1, m]$ ,  $\|\mathbf{a}_{*i}\| = \sqrt{\sum_{p=1}^n a_{pi}^2} \leq \sqrt{n}\rho\varepsilon$  as  $|a_{pi}| \leq \rho\varepsilon$ . Therefore,  $\|\mathbf{A}_T\| \leq \sqrt{mn}\rho\varepsilon$ . Then,  $\left\| \frac{1}{\lambda_{m_0+1}^2} \mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T \right\| \leq \|\mathbf{A}_T\|^2 \|\mathbf{L}_1\|^2 / \lambda_{m_0+1}^2 = \|\mathbf{A}_T\|^2 / \lambda_{m_0+1}^2 = O(n\rho^2\varepsilon^2 / \lambda_{m_0+1}^2)$ .

By Lemma 6.5, it holds with probability  $1 - O(n^{-1/3})$  that the diagonal entries of  $\widehat{\Lambda}_0^\dagger$  are  $\Theta(1/(\sqrt{n}\rho^2\varepsilon^2))$ . Thus, it holds with probability  $1 - O(n^{-1/3})$  that  $\|\mathbf{A}_T^t \mathbf{L}_0 \left( \widehat{\Lambda}_0^\dagger \right)^2 \mathbf{L}_0^t \mathbf{A}_T\| = \|\widehat{\Lambda}_0^\dagger \mathbf{L}_0^t \mathbf{A}_T\|^2 \leq \|\widehat{\Lambda}_0^\dagger\|^2 \|\mathbf{L}_0\|^2 \|\mathbf{A}_T\|^2 = \|\mathbf{A}_T\|^2 / \Theta(n\rho^4\varepsilon^4) = O(1/(\rho^2\varepsilon^2))$ , which is  $O(n\rho^2\varepsilon^4 / \lambda_{m_0+1}^2)$  as  $\lambda_{m_0+1} = O(\sqrt{n}\rho^2\varepsilon^3)$  with probability  $1 - O(n^{-1/3})$  by Lemma 6.6.

We conclude from the previous two paragraphs that  $\|\mathbf{H}_T \mathbf{H}_T^t\| = O(n\rho^2\varepsilon^2/\lambda_{m_0+1}^2)$  or, equivalently,  $\|\mathbf{H}_T\| = O(\sqrt{n}\rho\varepsilon/\lambda_{m_0+1})$ . In the rest of the proof, we show that every eigenvalue of  $\mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T$  is  $\Omega(n\rho^2\varepsilon^2)$ . This implies that the minimum eigenvalue of  $\mathbf{H}_T \mathbf{H}_T^t$  is  $\Omega(n\rho^2\varepsilon^2/\lambda_{m_0+1}^2)$  or, equivalently, the minimum singular value of  $\mathbf{H}_T$  is  $\Omega(\sqrt{n}\rho\varepsilon/\lambda_{m_0+1})$ .

Lemmas 7.7 and 7.8 imply that, with probability  $1 - O(n^{-1/3})$ , every column vector  $\mathbf{a}_{*i}$  of  $\mathbf{A}_T$  makes an angle at least  $\arccos(\frac{1}{9}m^{-1}n^{-1/3}) - O(\varepsilon^2)$  with the column space of  $\mathbf{L}_0$ . Therefore, with probability  $1 - O(n^{-1/3})$ , every column vector  $\mathbf{a}_{*i}$  of  $\mathbf{A}_T$  makes an angle at most  $O(\varepsilon^2) + \arcsin(\frac{1}{9}m^{-1}n^{-1/3})$  with the column space of  $\mathbf{L}_1$ . Let  $\bar{\mathbf{a}}_{*i}$  denote the projection of  $\mathbf{a}_{*i}$  in the column space of  $\mathbf{L}_1$ .

Take an arbitrary pair of distinct columns  $\mathbf{a}_{*i}$  and  $\mathbf{a}_{*j}$  of  $\mathbf{A}_T$ . By Lemma 7.6, the acute angle between the support lines of  $\mathbf{a}_{*i}$  and  $\mathbf{a}_{*j}$  is at least  $\arccos(\frac{1}{9}m^{-1}n^{-1/3})$ . We have argued in the previous paragraph that each of  $\mathbf{a}_{*i}$  and  $\mathbf{a}_{*j}$  makes an angle at most  $O(\varepsilon^2) + \arcsin(\frac{1}{9}m^{-1}n^{-1/3})$  with the column space of  $\mathbf{L}_1$ . Therefore, the acute angle between the support lines of  $\bar{\mathbf{a}}_{*i}$  and  $\bar{\mathbf{a}}_{*j}$  is at least  $\pi/2 - O(\varepsilon^2) - 3\arcsin(\frac{1}{9}m^{-1}n^{-1/3})$ .

For every  $i \in [1, m]$ ,  $\bar{\mathbf{a}}_{*i} = \mathbf{L}_1 \mathbf{L}_1^t \mathbf{a}_{*i}$ . Thus,  $(\bar{\mathbf{a}}_{*1} \ \cdots \ \bar{\mathbf{a}}_{*m}) = \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T$ . It follows that:

$$\begin{aligned} \mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T &= (\mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T)^t \cdot (\mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T) \\ &= \begin{pmatrix} \sum_{p=1}^n \bar{a}_{p1}^2 & \sum_{p=1}^n \bar{a}_{p1} \bar{a}_{p2} & \cdots & \sum_{p=1}^n \bar{a}_{p1} \bar{a}_{pm} \\ \sum_{p=1}^n \bar{a}_{p2} \bar{a}_{p1} & \sum_{p=1}^n \bar{a}_{p2}^2 & \cdots & \sum_{p=1}^n \bar{a}_{p2} \bar{a}_{pm} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{p=1}^n \bar{a}_{pm} \bar{a}_{p1} & \sum_{p=1}^n \bar{a}_{pm} \bar{a}_{p2} & \cdots & \sum_{p=1}^n \bar{a}_{pm}^2 \end{pmatrix} \end{aligned}$$

By the Gershgorin Circle Theorem, the minimum eigenvalue of  $\mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T$  is greater than or equal to  $\min_{i \in [1, m]} \left\{ \sum_{p=1}^n \bar{a}_{pi}^2 - \sum_{j \in [1, m] \setminus \{i\}} |\sum_{p=1}^n \bar{a}_{pi} \bar{a}_{pj}| \right\}$ .

We have shown earlier that  $\|\bar{\mathbf{a}}_{*i}\| \geq \|\mathbf{a}_{*i}\| \cos(O(\varepsilon^2) + \arcsin(\frac{1}{9}m^{-1}n^{-1/3}))$ , which is at least  $\|\mathbf{a}_{*i}\| (1 - \frac{1}{8}m^{-1}n^{-1/3})$  for a small enough  $\varepsilon$ . Note that  $\sum_{p=1}^n \bar{a}_{pi}^2 = \|\bar{\mathbf{a}}_{*i}\|^2$ . We have also shown that  $|\bar{\mathbf{a}}_{*i}^t \cdot \bar{\mathbf{a}}_{*j}| \leq \|\bar{\mathbf{a}}_{*i}\| \|\bar{\mathbf{a}}_{*j}\| \cos(\pi/2 - O(\varepsilon^2) - 3\arcsin(\frac{1}{9}m^{-1}n^{-1/3}))$ , which is at most  $\|\bar{\mathbf{a}}_{*i}\| \|\bar{\mathbf{a}}_{*j}\| \cdot \frac{1}{2}m^{-1}n^{-1/3} \leq \|\mathbf{a}_{*i}\| \|\mathbf{a}_{*j}\| \cdot \frac{1}{2}m^{-1}n^{-1/3}$  for a small enough  $\varepsilon$ . Note that  $\sum_{p=1}^n \bar{a}_{pi} \bar{a}_{pj} = \bar{\mathbf{a}}_{*i}^t \cdot \bar{\mathbf{a}}_{*j}$ .

Apply Lemma 5.4(iii) with  $c = \frac{1}{10(m+2)}$  to show that the relation  $\left| \|\mathbf{a}_{*i}\|^2 - \frac{1}{m+2}n\rho^2\varepsilon^2 \right| \leq \frac{1}{10(m+2)}n^{2/3}\rho^2\varepsilon^2 + O(n\rho^2\varepsilon^4)$  holds with probability  $1 - O(n^{-1/3})$ . Thus,  $\sum_{p=1}^n \bar{a}_{pi}^2 = \|\bar{\mathbf{a}}_{*i}\|^2 \geq \|\mathbf{a}_{*i}\|^2 (1 - \frac{1}{4}n^{-1/3}) \geq \frac{1}{m+2}n\rho^2\varepsilon^2 - \frac{1}{10(m+2)}n^{2/3}\rho^2\varepsilon^2 - \frac{1}{4(m+2)}n^{2/3}\rho^2\varepsilon^2 - \frac{1}{40(m+2)}n^{1/3}\rho^2\varepsilon^2 - O(n\rho^2\varepsilon^4)$ . Also,  $\sum_{j \in [1, m] \setminus \{i\}} |\sum_{p=1}^n \bar{a}_{pi} \bar{a}_{pj}| = \sum_{j \in [1, m] \setminus \{i\}} |\mathbf{a}_{*i}^t \cdot \mathbf{a}_{*j}| \leq \sum_{j \in [1, m] \setminus \{i\}} \|\mathbf{a}_{*i}\| \cdot \|\mathbf{a}_{*j}\| \cdot \frac{1}{2}m^{-1}n^{-1/3} \leq (m-1) \cdot \left( \frac{11}{10(m+2)}n\rho^2\varepsilon^2 + O(n\rho^2\varepsilon^4) \right) \cdot \frac{1}{2}m^{-1}n^{-1/3} < \frac{11}{20(m+2)}n^{2/3}\rho^2\varepsilon^2 + O(n^{2/3}\rho^2\varepsilon^4)$ . Hence, with probability  $1 - O(n^{-1/3})$ , the minimum eigenvalue of  $\mathbf{A}_T^t \mathbf{L}_1 \mathbf{L}_1^t \mathbf{A}_T$  is greater than or equal to  $\frac{1}{m+2}n\rho^2\varepsilon^2 - \frac{1}{m+2} \left( \frac{1}{10} + \frac{1}{4} + \frac{1}{40} + \frac{11}{20} \right) n\rho^2\varepsilon^2 - O(n\rho^2\varepsilon^4) = \frac{3}{40}n\rho^2\varepsilon^2 - O(n\rho^2\varepsilon^4)$  which is  $\Omega(n\rho^2\varepsilon^2)$  for a small enough  $\varepsilon$ .  $\square$

We are now ready to bound the eigenvalues of  $\mathbf{H}\mathbf{H}^t$ —the main result of this section.

**Lemma 7.10** *If  $\varepsilon$  is sufficiently small, then with probability  $1 - O(n^{-1/3})$ , the  $m$  largest eigenvalues of  $\mathbf{H}\mathbf{H}^t$  are  $\Theta(n\rho^2\varepsilon^2/\lambda_{m_0+1}^2)$  and the  $d-m$  smallest eigenvalues of  $\mathbf{H}\mathbf{H}^t$  are  $O(n\rho^2\varepsilon^4/\lambda_{m_0+1}^2)$ .*

*Proof.* By Lemma 7.4, we can rotate  $\mathbb{R}^d$  so that the coordinate axes  $x_1, \dots, x_m$  span  $\mathcal{T}$ . Partition  $\mathbf{H}\mathbf{H}^t$  into blocks using  $\mathbf{H}_T$  and  $\mathbf{H}_N$ :

$$\mathbf{H}\mathbf{H}^t = \begin{pmatrix} \mathbf{H}_T \mathbf{H}_T^t & \mathbf{H}_T \mathbf{H}_N^t \\ \mathbf{H}_N \mathbf{H}_T^t & \mathbf{H}_N \mathbf{H}_N^t \end{pmatrix}$$



We apply the generalization of the Gershgorin circle theorem [13]. Let  $G_1$  be the set of all real numbers  $\mu$  such that  $(\|(\mathbf{H}_T \mathbf{H}_T^t - \mu \mathbf{I}_m)^{-1}\|)^{-1} \leq \|\mathbf{H}_T \mathbf{H}_N^t\|$ , and let  $G_2$  be the set of all real numbers  $\mu$  such that  $(\|(\mathbf{H}_N \mathbf{H}_N^t - \mu \mathbf{I}_{d-m})^{-1}\|)^{-1} \leq \|\mathbf{H}_N \mathbf{H}_T^t\|$ .

By (6.2) and (6.3), the numbers in  $G_1$  are at least the minimum eigenvalue value of  $\mathbf{H}_T \mathbf{H}_T^t$  minus  $\|\mathbf{H}_T \mathbf{H}_N^t\|$  and at most the maximum eigenvalue value of  $\mathbf{H}_T \mathbf{H}_T^t$  plus  $\|\mathbf{H}_T \mathbf{H}_N^t\|$ . If  $\varepsilon$  is small enough, the numbers in  $G_1$  are  $\Theta(n \varrho^2 \varepsilon^2 / \lambda_{m_0+1}^2)$  with probability  $1 - O(n^{-1/3})$  because Lemmas 7.5 and 7.9 imply that, with probability  $1 - O(n^{-1/3})$ , the eigenvalues of  $\mathbf{H}_T \mathbf{H}_T^t$  are  $\Theta(n \varrho^2 \varepsilon^2 / \lambda_{m_0+1}^2)$ , which dominates  $\|\mathbf{H}_T \mathbf{H}_N^t\| \leq \|\mathbf{H}_T\| \|\mathbf{H}_N\| = O(n \varrho^2 \varepsilon^4 / \lambda_{m_0+1}^2)$ .

By (6.3), the numbers in  $G_2$  are at most the maximum eigenvalue value of  $\mathbf{H}_N \mathbf{H}_N^t$  plus  $\|\mathbf{H}_N \mathbf{H}_T^t\|$ . If  $\varepsilon$  is small enough, the numbers in  $G_2$  are  $O(n \varrho^2 \varepsilon^4 / \lambda_{m_0+1}^2)$  with probability  $1 - O(n^{-1/3})$  because Lemmas 7.5 and 7.9 imply that, with probability  $1 - O(n^{-1/3})$ ,  $\|\mathbf{H}_N \mathbf{H}_N^t\| = \|\mathbf{H}_N\|^2 = O(n \varrho^2 \varepsilon^6 / \lambda_{m_0+1}^2)$  and  $\|\mathbf{H}_N \mathbf{H}_T^t\| \leq \|\mathbf{H}_N\| \|\mathbf{H}_T\| = O(n \varrho^2 \varepsilon^4 / \lambda_{m_0+1}^2)$ .

As a result, if  $\varepsilon$  is small enough, then with probability  $1 - O(n^{-1/3})$ , all numbers in  $G_2$  are smaller than those in  $G_1$ . By Lemma 6.1, the disjointness of  $G_1$  and  $G_2$  implies that  $G_1$  and  $G_2$  contain exactly  $m$  and  $d-m$  eigenvalues of  $\mathbf{H} \mathbf{H}^t$ , respectively. Hence, the lemma follows.  $\square$

## 8 Proof of Lemma 2.1

Let  $S$  be a set of sample points in  $\mathcal{M}$ . Without loss of generality, we assume that the origin is a sample point in  $S$ . We require  $S$  to satisfy the condition that the distance between the origin and its  $(n+1)$ -th nearest sample point is at most  $\varrho \varepsilon$  for some sufficiently small  $\varepsilon \in (0, \varepsilon_0]$  where  $\varepsilon_0$  is the constant in Lemma 5.1. We rotate  $\mathbb{R}^d$  so that the coordinate axes  $x_1, \dots, x_m$  span the tangent space  $\mathcal{T}$  of  $\mathcal{M}$  at the origin. By Lemma 6.2, this does not change the singular values of  $\mathbf{B}$ .

Let  $E$  denote the event that  $\lambda_{m_0+1} = 0$  and there exists  $q \in [1, n]$  such that  $a_{q,k} \neq 0$  for some  $k \in [m+1, d]$ . Our goal is to prove that  $\Pr(E) = O(n^{-1/3})$ . Then, it holds with probability  $1 - O(n^{-1/3})$  that if  $\lambda_{m_0+1} = 0$ , we have  $a_{p,i} = 0$  for all  $p \in [1, n]$  and for all  $i \in [m+1, d]$ , and so Lemma 2.1 is true.

By Lemma 6.3, it holds with probability  $1 - O(n^{-1/3})$  that the eigenvalues of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  are at least  $c_0 n \varrho^4 \varepsilon^4$  and at most  $c_1 n \varrho^4 \varepsilon^4$  for some constants  $c_0$  and  $c_1$ . The proof of Lemma 6.3 reveals that  $c_0$  and  $c_1$  are polynomials in  $m$ . Let  $F$  denote the event that the eigenvalues of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  lie between  $c_0 n \varrho^4 \varepsilon^4$  and  $c_1 n \varrho^4 \varepsilon^4$ . The probability  $\Pr(E)$  can be split up into the following sum:

$$\Pr(E|F) \cdot \Pr(F) + \Pr(E|\neg F) \cdot \Pr(\neg F).$$

The second term is  $O(n^{-1/3})$  because  $\Pr(\neg F) = O(n^{-1/3})$  by Lemma 6.3. We show that  $\Pr(E|F) = 0$  below. From now on, we assume that the condition  $F$  holds.

When the event  $E$  happens,  $\lambda_{m_0+1} = 0$  and there exists index  $q \in [1, n]$  such that  $a_{q,k} \neq 0$  for some  $k \in [m+1, d]$ . By swapping coordinate axes if necessary, we can further assume that  $a_{q,m+1} \neq 0$ . Let  $\mathbf{C}_{q,TT}$  denote the matrix obtained by deleting the row in  $\mathbf{B}_{TT}$  for  $\mathbf{a}_q$ . Since  $n \geq m_0 + 1$ , there are at least  $m_0$  rows in  $\mathbf{C}_{q,TT}$ . We claim that  $\mathbf{C}_{q,TT}$  has rank  $m_0$ , provided that  $n > c_0 m_0$ . Otherwise, the smallest eigenvalue of  $\mathbf{C}_{q,TT}^t \mathbf{C}_{q,TT}$  is zero, which implies that there is a unit direction  $\mathbf{u} \in \mathbb{R}^{m_0}$  that is orthogonal to every row vector of  $\mathbf{C}_{q,TT}$ . Since every coordinate of  $\mathbf{a}_q$  has magnitude at most  $\varrho \varepsilon$ , the projection of the row for  $\mathbf{a}_q$  in  $\mathbf{B}_{TT}$  onto  $\mathbf{u}$  has a squared length at most  $m_0 \varrho^4 \varepsilon^4$ . Therefore, if we project the row vectors of  $\mathbf{B}_{TT}$  onto  $\mathbf{u}$ , the sum of the squared lengths of the projections is at most  $m_0 \varrho^4 \varepsilon^4$ , which implies that the smallest eigenvalue of  $\mathbf{B}_{TT}^t \mathbf{B}_{TT}$  is at most  $m_0 \varrho^4 \varepsilon^4 < c_0 n \varrho^4 \varepsilon^4$ , a contradiction to the condition  $F$ .

Let  $\mathbf{C}_q$  denote the matrix obtained from  $\mathbf{B}$  by removing the row for  $\mathbf{a}_q$ . Since  $\mathbf{C}_{q,TT}$  is a

submatrix of  $C_q$  which is in turn a submatrix of  $B$ , we have

$$m_0 = \text{rank}(C_{q,TT}) \leq \text{rank}(C_q) \leq \text{rank}(B).$$

Since  $\lambda_{m_0+1} = 0$ , the rank of  $B$  is at most  $m_0$ , and so  $\text{rank}(B) \geq \text{rank}(C_{q,TT}) = m_0$  implies that  $\text{rank}(B) = m_0$ . This allows us to conclude that  $C_q$  has rank  $m_0$ .

Since  $\text{rank}(B) = m_0$ , the column space of  $B$  has rank  $m_0$ , which implies that every column in  $B$  is a linear combination of the columns of  $B_{TT}$ . Therefore, there exist coefficients  $g_{ij}$ 's such that

$$\forall p \in [1, n], \quad \frac{1}{\sqrt{2}} a_{p,m+1}^2 = \sum_{i \leq j \in [1, m]} g_{ij} a_{pi} a_{pj}. \quad (8.1)$$

Since  $n \geq m_0 + 1$ , there are  $m_0$  coefficients  $g_{ij}$ 's, and  $C_q$  has rank  $m_0$ , the coefficients  $g_{ij}$ 's are completely determined by the following smaller system:

$$\forall p \in [1, n] \setminus \{q\}, \quad \frac{1}{\sqrt{2}} a_{p,m+1}^2 = \sum_{i \leq j \in [1, m]} g_{ij} a_{pi} a_{pj}. \quad (8.2)$$

In other words, when  $E$  happens under condition  $F$ , the coefficients  $g_{ij}$ 's are determined irrespective of the coordinates of  $\mathbf{a}_q$  and yet  $\mathbf{a}_q$  must satisfy (8.1). We can interpret (8.1) as saying that the sample points  $\mathbf{a}_p$ ,  $p \in [1, n]$ , lie in the following hypersurface:

$$H(\mathbf{x}) = \frac{1}{\sqrt{2}} x_{m+1}^2 - \sum_{i \leq j \in [1, m]} g_{ij} x_i x_j = 0. \quad (8.3)$$

Both the hypersurface  $H(\mathbf{x}) = 0$  and  $\mathcal{M}$  contain the origin. Let  $L$  denote the linear subspace spanned by axes  $x_1, \dots, x_{m+1}$ . The intersection of  $L$  and the hypersurface  $H(\mathbf{x}) = 0$  is a conic surface, and this can be seen as follows. Recall that the axes  $x_1, \dots, x_m$  span the true tangent space  $\mathcal{T}$  of  $\mathcal{M}$  at the origin. Take any unit vector  $\mathbf{u} \in \mathcal{T}$  that makes an angle  $\theta_i$  with the axis  $x_i$  for  $i \in [1, m]$ . For any  $c \in \mathbb{R}$ , the  $x_{m+1}$  coordinate of the point in the hypersurface  $H(\mathbf{x}) = 0$  that projects to the point  $c\mathbf{u} \in \mathcal{T}$  can be written as

$$\begin{aligned} x_{m+1}^2 &= \left( \sum_{i \leq j \in [1, m]} \sqrt{2} g_{ij} \cos \theta_i \cos \theta_j \right) c^2 \\ \Rightarrow x_{m+1} &= \pm \left( \sum_{i \leq j \in [1, m]} \sqrt{2} g_{ij} \cos \theta_i \cos \theta_j \right)^{1/2} c. \end{aligned} \quad (8.4)$$

Therefore, the cross-section of  $H(\mathbf{x}) = 0$  in the plane spanned by  $\mathbf{u}$  and the axis  $x_{m+1}$  consists of two lines through the origin with slopes of the same magnitude but opposite signs.

Recall the coordinate function  $f_{m+1} : \mathbb{R}^m \rightarrow \mathbb{R}$  for  $\mathcal{M}$  such that, given a point  $(x_1 \dots x_d)^t \in \mathcal{M}$ , we have  $x_{m+1} = f_{m+1}((x_1 \dots x_m)^t)$  in a local neighborhood of the origin. We can choose the constant  $\varepsilon_0$  in Lemma 5.1 so that  $\varrho\varepsilon_0$  is at most the radius of this local neighborhood. Also, recall that the Taylor expansion of  $f_{m+1}$  does not have a constant or a linear term, that is,

$$\forall \mathbf{u} \in \mathbb{R}^m, \forall c \in \mathbb{R}, \quad f_{m+1}(c\mathbf{u}) = \frac{c^2}{2} D^2 f_{m+1}|_0(\mathbf{u}, \mathbf{u}) + \dots \quad (8.5)$$

We claim that the set of points  $\mathcal{K} = \{\mathbf{x} \in \mathcal{M} : x_{m+1} \neq 0 \wedge H(\mathbf{x}) = 0\}$  has measure zero. Suppose not. Then, there exists a unit vector  $\mathbf{v} = (v_1 \dots v_m)^t \in \mathcal{T}$  such that the intersection between  $\mathcal{K}$  and the plane spanned by  $\mathbf{v}$  and the coordinate axis  $x_{m+1}$  consists of some curve

segment(s) of positive length. By the definition of  $\mathcal{K}$ ,  $f_{m+1}$  is not identically zero along the directions  $\mathbf{v}$  and  $-\mathbf{v}$ . Then, since (8.4) is a linear equation, it can only agree with (8.5) at isolated values of  $c$  after we substitute  $\mathbf{u}$  by  $\mathbf{v}$ . This is a contradiction to the existence of curve segment(s) of positive length in the intersection of  $\mathcal{K}$  and the plane spanned by  $\mathbf{v}$  and the axis  $x_{m+1}$ . We conclude that  $\mathcal{K}$  has measure zero. It follows that the probability of drawing the sample point  $\mathbf{a}_q$  from  $\mathcal{K}$  is zero, which implies that  $\Pr(E|F) = 0$ .

In summary,  $\Pr(E|F) \cdot \Pr(F) + \Pr(E|\neg F) \cdot \Pr(\neg F) = O(n^{-1/3})$ , completing the proof of Lemma 2.1.

## 9 Proof of Theorem 1.1

By Lemma 7.4, we can assume that  $\mathbb{R}^d$  has been rotated so that the coordinate axes  $x_1, \dots, x_m$  span the tangent space  $\mathcal{T}$  of  $\mathcal{M}$  at the origin. Then, we can partition  $\mathbf{H} = \begin{pmatrix} \mathbf{H}_T \\ \mathbf{H}_N \end{pmatrix}$ .

We call  $\text{TANGENT}(\mathbf{A})$ . If  $\lambda_{m_0+1} = 0$ , then Lemma 2.1 implies that, with probability  $1 - O(n^{-1/3})$ , the estimated tangent space is equal to  $\mathcal{T}$ , and therefore, there is no angular error.

Assume that  $\lambda_{m_0+1} > 0$  for the rest of the proof. Then, the estimated tangent space is spanned by the unit eigenvectors corresponding to the  $m$  largest eigenvalues of  $\mathbf{A}^t \mathbf{L} \Sigma \mathbf{L}^t \mathbf{A} = \mathbf{H} \mathbf{H}^t$ . Let  $\mathbf{e}$  be one of these  $m$  unit eigenvectors. Divide  $\mathbf{e}$  into two parts  $(\mathbf{v}^t \ \mathbf{w}^t)^t$ , where  $\mathbf{v}$  consists of the first  $m$  coordinates and  $\mathbf{w}$  consists of the last  $d - m$  coordinates. Let  $\sigma$  be the eigenvalue of  $\mathbf{H} \mathbf{H}^t$  corresponding to  $\mathbf{e}$ . Then,

$$\mathbf{H} \mathbf{H}^t \mathbf{e} = \begin{pmatrix} \mathbf{H}_T \mathbf{H}_T^t & \mathbf{H}_T \mathbf{H}_N^t \\ \mathbf{H}_N \mathbf{H}_T^t & \mathbf{H}_N \mathbf{H}_N^t \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix} = \sigma \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix} \quad (9.1)$$

We want to bound  $\arctan(\|\mathbf{w}\|/\|\mathbf{v}\|)$ , the angle between  $\mathbf{e}$  and  $\mathcal{T}$ .

We first show that  $\mathbf{v} \neq \mathbf{0}_{m,1}$  with probability  $1 - O(n^{-1/3})$ . Suppose that  $\mathbf{v} = \mathbf{0}_{m,1}$ . By (9.1), we get  $\mathbf{H}_N \mathbf{H}_N^t \mathbf{w} = \sigma \mathbf{w}$ , meaning that  $\sigma$  is also an eigenvalue of  $\mathbf{H}_N \mathbf{H}_N^t$  in addition to being one of the  $m$  largest eigenvalues of  $\mathbf{H} \mathbf{H}^t$ . This occurs with probability  $O(n^{-1/3})$  because, with probability  $1 - O(n^{-1/3})$ , the  $m$  largest eigenvalues of  $\mathbf{H} \mathbf{H}^t$  are  $\Theta(n \varrho^2 \varepsilon^2 / \lambda_{m_0+1}^2)$  by Lemma 7.10, but the eigenvalues of  $\mathbf{H}_N \mathbf{H}_N^t$  are  $O(n \varrho^2 \varepsilon^6 / \lambda_{m_0+1}^2)$  by Lemma 7.5.

Assume that  $\mathbf{v}$  is non-zero. By (9.1),  $\mathbf{w} = (\sigma \mathbf{I}_{d-m} - \mathbf{H}_N \mathbf{H}_N^t)^{-1} \mathbf{H}_N \mathbf{H}_T^t \mathbf{v}$ . By Lemmas 7.5 and 7.9, it holds with probability  $1 - O(n^{-1/3})$  that  $\|\mathbf{H}_N \mathbf{H}_T^t\| \leq \|\mathbf{H}_N\| \|\mathbf{H}_T\| = O(n \varrho^2 \varepsilon^4 / \lambda_{m_0+1}^2)$ . By Lemmas 7.5 and 7.10, it holds with probability  $1 - O(n^{-1/3})$  that the eigenvalues of  $\mathbf{H}_N \mathbf{H}_N^t$  are  $O(n \varrho^2 \varepsilon^6 / \lambda_{m_0+1}^2)$  and  $\sigma = \Theta(n \varrho^2 \varepsilon^2 / \lambda_{m_0+1}^2)$ . The smallest eigenvalue of  $\sigma \mathbf{I}_{d-m} - \mathbf{H}_N \mathbf{H}_N^t$  is thus  $\Theta(n \varrho^2 \varepsilon^2 / \lambda_{m_0+1}^2)$ , implying that  $\|(\sigma \mathbf{I}_{d-m} - \mathbf{H}_N \mathbf{H}_N^t)^{-1}\| = O(\lambda_{m_0+1}^2 / (n \varrho^2 \varepsilon^2))$ . Hence,  $\|\mathbf{w}\| \leq \|(\sigma \mathbf{I}_{d-m} - \mathbf{H}_N \mathbf{H}_N^t)^{-1}\| \cdot \|\mathbf{H}_N \mathbf{H}_T^t\| \cdot \|\mathbf{v}\| = O(\varepsilon^2) \cdot \|\mathbf{v}\|$ . The angle between  $\mathbf{e}$  and  $\mathcal{T}$  is  $\arctan(\|\mathbf{w}\|/\|\mathbf{v}\|) = O(\varepsilon^2)$ .

We conclude that, with probability at least  $1 - m \cdot O(n^{-1/3}) = 1 - O(n^{-1/3})$ , all eigenvectors corresponding to the  $m$  largest eigenvalues of  $\mathbf{H} \mathbf{H}^t$  make an  $O(\varepsilon^2)$  angle with  $\mathcal{T}$ . Let  $(\mathbf{U} \ \mathbf{V})$  be a  $d \times d$  orthonormal matrix such that the columns of  $\mathbf{V}$  form an orthonormal basis of  $\mathcal{T}$ . This implies that, with probability  $1 - O(n^{-1/3})$ , all eigenvectors corresponding to the  $m$  largest eigenvalues of  $\mathbf{H} \mathbf{H}^t$  make an  $O(\varepsilon^2)$  angle with the column space of  $\mathbf{V}$  and hence an  $\pi/2 - O(\varepsilon^2)$  angle with the column space of  $\mathbf{U}$ . Then, Lemma 7.2 implies that  $\sqrt{m} \cdot O(\varepsilon^2) = O(\varepsilon^2)$  is an upper bound on the angle between  $\mathcal{T}$  and the space spanned by the eigenvectors corresponding to the  $m$  largest eigenvalues of  $\mathbf{H} \mathbf{H}^t$ .

## 10 Conclusion

We present an algorithm for estimating tangent spaces from a given set of sample points in an unknown manifold. The algorithm works locally and uses the  $n$  sample points nearest to  $\mathbf{p}$ . The distance from  $\mathbf{p}$  to the  $(n + 1)$ -th nearest sample point can be expressed as  $\varrho\varepsilon$ , where  $\varepsilon \in (0, 1)$  and  $\varrho$  is the local feature size at  $\mathbf{p}$ . (The algorithm does not need to know  $\varrho$  though.) When we fix  $n$ , the value  $\varepsilon$  decreases as sampling density increases. Assuming that the sample points are distributed according to a Poisson process with an unknown parameter, our algorithm guarantees an  $O(\varepsilon^2)$  bound on the angular error with high probability. The quadratic angular error convergence has been confirmed in our experiments.

The angular error bounds in [2, 10, 11, 15, 26] hold for all sample points. Our  $O(\varepsilon^2)$  angular error bound applies to the center of a local neighborhood in which the sample points are used to estimate the tangent space at the center. One should be able to extend our result so that it applies simultaneously to centers of disjoint neighborhoods by restricting the range of  $\varepsilon$  further and estimating the failure probability using the union bound. Further work is needed to see if our angular error bound can be guaranteed at all sample points.

Our algorithm assumes that the manifold dimension  $m$  is known to us. The results in Section 6 show that the largest  $\binom{m+1}{2}$  eigenvalues of  $\mathbf{B}^t\mathbf{B}$  are  $\Theta(n\varrho^4\varepsilon^4)$  and the other eigenvalues are  $O(n\varrho^4\varepsilon^6)$ . Therefore, one should be able to determine the manifold dimension automatically by detecting this  $\Omega(\varepsilon^2)$  factor gap in the eigenvalues. This approach of finding gaps in the spectrum of eigenvalues has been used in several previous work on detecting manifold dimension.

One may wonder what happens if we omit the conversion of  $\mathbf{B}$  to  $\widehat{\mathbf{B}}$  in practice. This is equivalent to defining  $\Sigma$  as  $\text{diag}_n\left(\frac{1}{\lambda_1^2}, \frac{1}{\lambda_2^2}, \dots, \frac{1}{\lambda_n^2}\right)$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  are the singular values of  $\mathbf{B}$ . We experimented with this alternative method. While this method performs reasonably in the noiseless case, it fails badly in the noisy cases. So the conversion of  $\mathbf{B}$  to  $\widehat{\mathbf{B}}$  makes a real difference. There is also the possibility that  $\lambda_{m_0+1}$  is positive but close to zero. In that case, although the theoretical analysis holds, there will be numerical issues in forming  $\Sigma = \text{diag}_n\left(\frac{1}{\lambda_1^2}, \dots, \frac{1}{\lambda_{m_0}^2}, \frac{1}{\lambda_{m_0+1}^2}, \dots, \frac{1}{\lambda_{m_0+1}^2}\right)$ . Therefore, some thresholding of the singular values may be necessary for numerical stability.

## Acknowledgment

We thank the anonymous referees for their helpful comments and useful advice.

## References

- [1] M. Belkin and P. Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. *Advances in Neural Information Processing Systems*, 14 (2002).
- [2] M. Belkin, J. Sun, and Y. Wang. Constructing Laplace operator from point clouds in  $R^d$ . *Proceedings of the 20th Annual ACM-SIAM Symposium on Discrete Algorithms*, 2009, 1021–1040.
- [3] J.-D. Boissonnat and A. Ghosh. Manifold reconstruction using tangential Delaunay complexes. *Proceedings of the 26th Annual Symposium on Computational Geometry*, 2010, 324–333.
- [4] J.-D. Boissonnat, L.J. Guibas and S.Y. Oudot. Manifold reconstruction in arbitrary dimensions using witness complexes. *Discrete and Computational Geometry*, 42 (2009), 37–70.

- [5] O.J. Boxma and U. Yechiali. Poisson Processes, Ordinary and Compound. *Encyclopedia of Statistics in Quality and Reliability*, F. Ruggeri, R.S. Kenett and F.W. Faltin (eds.), Wiley, New York, 2007.
- [6] K.M. Carter, R. Raich, and A.O. Hero. On local intrinsic dimension estimation and its applications. *IEEE Transactions on Signal Processing*, 58 (2010), 650–663.
- [7] T.F. Chan. An improved algorithm for computing the singular value decomposition. *ACM Transactions on Mathematical Software*, 8 (1982), 72–83.
- [8] S.-W. Cheng and M-K. Chiu. Dimension detection via slivers. *Proceedings of the 20th Annual ACM-SIAM Symposium on Discrete Algorithms*, 2009, 1001-1010.
- [9] S.-W. Cheng, T.K. Dey and E.A. Ramos. Manifold reconstruction from point samples. *Proceedings of the 16th Annual ACM-SIAM Symposium on Discrete Algorithms*, 2005, 1018–1027.
- [10] S.-W. Cheng, Y. Wang, and Z. Wu. Provable Dimension Detection using Principal Component Analysis. *International Journal of Computational Geometry and Applications*, 18 (2008), 414–440.
- [11] T. K. Dey, J. Giesen, S. Goswami and W. Zhao. Shape dimension and approximation from samples. *Discrete and Computational Geometry*, 29 (2003), 419–434.
- [12] S.C. Eisenstat and I.C.F. Ipsen. Relative perturbation bounds for eigenspaces and singular vector subspaces. *Proceedings of the 5th SIAM Conference on Applied Linear Algebra*, 1994, 62–66.
- [13] D.G. Feingold, R.S. Varga, Block diagonally dominant matrices and generalizations of the Gershgorin circle theorem, *Pacific Journal of Mathematics*, 12 (1962), 1241–1250.
- [14] M. Gashler and T. Martinez. Tangent space guided intelligent neighbor finding. *Proceedings of International Joint Conference on Neural Networks*, 2011, 2617–2624.
- [15] J. Giesen and U. Wagner. Shape dimension and intrinsic metric from samples of manifolds with high codimension. *Discrete and Computational Geometry*, 32 (2004), 245–267.
- [16] G.H. Golub and C. Reinsch. Singular value decomposition and least square solutions. In *Handbook for Automatic Computation*, II, *Linear Algebra*, J.H. Wilkinson and C. Reinsch (eds.), Springer-Verlag, New York, 1971.
- [17] G.H. Golub and C.F. van Loan. *Matrix Computations*, Johns Hopkins University Press, 1996.
- [18] D. Gong, X. Zhao, and G. Medioni. Robust multiple manifolds structure learning. *Proceedings of the International Conference on Machine Learning*, 2012,
- [19] I.S. Gradshteyn and I.M. Ryzhik. *Table of Integrals, Series, and Products*, Academic Press, 1994.
- [20] P.S. Heckbert and M. Garland. Optimal triangulation and quadric-based surface simplification. *Computational Geometry: Theory and Applications*, 14 (1999), 49–65.
- [21] M. Hein and J.-Y. Audibert. Intrinsic dimensionality estimation of submanifolds in Euclidean space. *Proceedings of the 22nd International Conference on Machine Learning*, 2005, 289–296.

- [22] B. Kégl. Intrinsic dimension estimation using packing numbers. *Advances in Neural Information Processing Systems*, 14 (2003), 681–688.
- [23] S. Lang. *Calculus of Several Variables*, Springer-Verlag, 1987.
- [24] F. Le Gall. Faster Algorithms for Rectangular Matrix Multiplication. *Proceedings of the 53rd IEEE Annual Symposium on Foundations of Computer Science*, 2012, 514–523.
- [25] E. Levina and P.J. Bickel. Maximum likelihood estimation of intrinsic dimension. *Advances in Neural Information Processing Systems*, 17 (2005), 777–784.
- [26] A.V. Little, M. Maggioni, and L. Rosasco. Multiscale geometric methods for data sets I: multiscale SVD, noise and curvature. *Computer Science and Artificial Intelligence Laboratory Technical Report*, MIT-CSAIL-TR-2012-029, CBCL-310, September 8, 2012.
- [27] A. massoud Farahmand, C. Szepesvári and J.-Y. Audibert. Manifold-adaptive dimension estimation. *Proceedings of the 24th International Conference on Machine Learning*, 2007, 265–272.
- [28] F. Morgan. *Riemannian Geometry: a beginner’s guide*, A. K. Peters, 1998.
- [29] J.C. Nascimento and J.G. Silva. Manifold learning for object tracking with multiple motion dynamics. *Proceedings of ECCV 2010*, Part III, LNCS 6313, 172–185.
- [30] S. Roweis and L. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290 (2000), 2323–2326.
- [31] B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10 (1998), 1299–1319.
- [32] G. Taubin. Estimation of planar curves, surfaces and nonplanar space curves defined by implicit equations, with applications to edge and range image segmentation, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13 (1991), 1115–1138.
- [33] J.B. Tenenbaum, V. de Silva and J.C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290 (2000), 2319–2323.