

# Spectral Clustering Based on Local PCA

**Ery Arias-Castro**

EARIASCA@UCSD.EDU

*Department of Mathematics  
University of California, San Diego  
La Jolla, CA 92093, USA*

**Gilad Lerman**

LERMAN@UMN.EDU

*Department of Mathematics  
University of Minnesota, Twin Cities  
Minneapolis, MN 55455, USA*

**Teng Zhang**

TENG.ZHANG@UCF.EDU

*Department of Mathematics  
University of Central Florida  
Orlando, FL 32816, USA*

**Editor:** Mikhail Belkin

## Abstract

We propose a spectral clustering method based on local principal components analysis (PCA). After performing local PCA in selected neighborhoods, the algorithm builds a nearest neighbor graph weighted according to a discrepancy between the principal subspaces in the neighborhoods, and then applies spectral clustering. As opposed to standard spectral methods based solely on pairwise distances between points, our algorithm is able to resolve intersections. We establish theoretical guarantees for simpler variants within a prototypical mathematical framework for multi-manifold clustering, and evaluate our algorithm on various simulated data sets.

**Keywords:** multi-manifold clustering, spectral clustering, local principal component analysis, intersecting clusters

## 1. Introduction

The task of multi-manifold clustering, where the data are assumed to be located near surfaces embedded in Euclidean space, is relevant in a variety of applications. In cosmology, it arises as the extraction of galaxy clusters in the form of filaments (curves) and walls (surfaces) (Valdarnini, 2001; Martínez and Saar, 2002); in motion segmentation, moving objects tracked along different views form affine or algebraic surfaces (Ma et al., 2008; Fu et al., 2005; Vidal and Ma, 2006; Chen et al., 2009); this is also true in face recognition, in the context of images of faces in fixed pose under varying illumination conditions (Ho et al., 2003; Basri and Jacobs, 2003; Epstein et al., 1995).

We consider a stylized setting where the underlying surfaces are nonparametric in nature, with a particular emphasis on situations where the surfaces intersect. Specifically, we assume the surfaces are smooth, for otherwise the notion of continuation is potentially ill-posed. For example, without smoothness assumptions, an L-shaped cluster is indistinguishable from the union of two line-segments meeting at right angle.

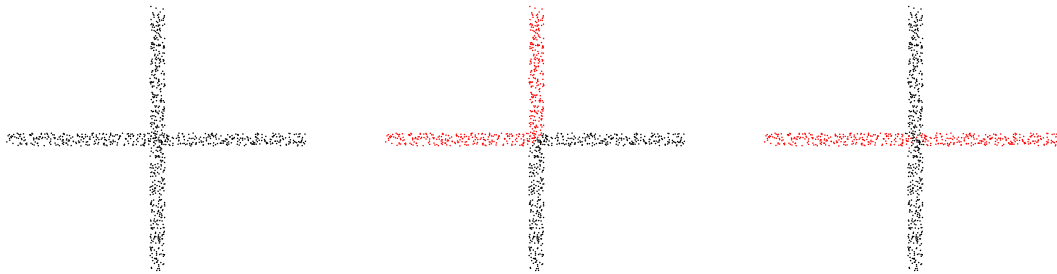


Figure 1: Two rectangular clusters intersecting at right angle. Left: the original data. Center: a typical output of the standard spectral clustering method of Ng et al. (2002), which is generally unable to resolve intersections. Right: a typical output of our method.

Spectral methods (Luxburg, 2007) are particularly suited for nonparametric settings, where the underlying clusters are usually far from convex, making standard methods like K-means irrelevant. However, a drawback of standard spectral approaches such as the well-known variant of Ng, Jordan, and Weiss (2002) is their inability to separate intersecting clusters. Indeed, consider the simplest situation where two straight clusters intersect at right angle, pictured in Figure 1 below. The algorithm of Ng et al. (2002) is based on pairwise affinities that are decreasing in the distances between data points, making it insensitive to smoothness and, therefore, intersections. And indeed, this algorithm typically fails to separate intersecting clusters, even in the easiest setting of Figure 1.

As argued in (Agarwal et al., 2005, 2006; Shashua et al., 2006), a multiway affinity is needed to capture complex structure in data (here, smoothness) beyond proximity attributes. For example, Chen and Lerman (2009a) use a flatness affinity in the context of *hybrid linear modeling*, where the surfaces are assumed to be affine subspaces, and subsequently extended to algebraic surfaces via the ‘kernel trick’ (Chen, Atev, and Lerman, 2009). Moving beyond parametric models, Arias-Castro, Chen, and Lerman (2011) introduce a localized measure of flatness.

Continuing this line of work, we suggest a spectral clustering method based on the estimation of the local linear structure (tangent bundle) via local principal component analysis (PCA). The idea of using local PCA combined with spectral clustering has precedents in the literature. In particular, our method is inspired by the work of Goldberg, Zhu, Singh, Xu, and Nowak (2009), where the authors develop a spectral clustering method within a semi-supervised learning framework. This approach is in the zeitgeist. While writing this paper, we became aware of two concurrent publications, by Wang, Jiang, Wu, and Zhou (2011) and by Gong, Zhao, and Medioni (2012), both proposing approaches very similar to ours.<sup>1</sup> We also mention the multiscale, spectral-flavored algorithm of Kushnir, Galun, and Brandt (2006), which is also based on local PCA. We comment on these spectral methods

---

1. The present paper was indeed developed in parallel with these two, first as a short version submitted to ICML in 2011.

in more detail later on. In fact, an early proposal also based on local PCA appears in the literature on subspace clustering (Fan et al., 2006)—although the need for localization is perhaps less intuitive in this setting.

The basic proposition of local PCA combined with spectral clustering has two main stages. The first one forms an affinity between a pair of data points that takes into account both their Euclidean distance and a measure of discrepancy between their tangent spaces. Each tangent space is estimated by PCA in a local neighborhood around each point. The second stage applies standard spectral clustering with this affinity. As a reality check, this relatively simple algorithm succeeds at separating the straight clusters in Figure 1. We tested our algorithm in more elaborate settings, some of them described in Section 4.

Other methods with a spectral component include those of Polito and Perona (2001) and Goh and Vidal (2007), which (roughly speaking) embed the points by a variant of LLE (Saul and Roweis, 2003) and then group the points by K-means clustering. There is also the method of Elhamifar and Vidal (2011), which chooses the neighborhood of each point by computing a sparse linear combination of the remaining points followed by an application of the spectral graph partitioning algorithm of Ng et al. (2002). Note that these methods work under the assumption that the surfaces do not intersect.

Besides spectral-type approaches to multi-manifold clustering, other methods appear in the literature. Among these, Gionis et al. (2005) and Haro et al. (2007) allow for intersecting surfaces but assume that they have different intrinsic dimension or density—and the proposed methodology is entirely based on such assumptions. We also mention the K-manifold method of Souvenir and Pless (2005), which propose an EM-type algorithm; and that of Guo et al. (2007), which propose to minimize an energy functional based on pairwise distances and local curvatures, leading to a combinatorial optimization.

Our contribution is the design and detailed study of a prototypical spectral clustering algorithm based on local PCA, tailored to settings where the underlying clusters come from sampling in the vicinity of smooth surfaces that may intersect. We endeavored to simplify the algorithm as much as possible without sacrificing performance. We provide theoretical results for simpler variants within a standard mathematical framework for multi-manifold clustering. To our knowledge, these are the first mathematically backed successes at the task of resolving intersections in the context of multi-manifold clustering, with the exception of (Arias-Castro et al., 2011), where the corresponding algorithm is shown to succeed at separating intersecting curves. The salient features of our algorithm are illustrated via numerical experiments.

The rest of the paper is organized as follows. In Section 2, we introduce our method in various variants. In Section 3, we analyze the simpler variants in a standard mathematical framework for multi-manifold learning. In Section 4, we perform some numerical experiments illustrating several features of our approach. In Section 5, we discuss possible extensions.

## 2. The Methodology

We introduce our algorithm and simpler variants that are later analyzed in a mathematical framework. We start with some review of the literature, zooming in on the most closely related publications.

## 2.1 Some Precedents

Using local PCA within a spectral clustering algorithm was implemented in four other publications we know of (Goldberg et al., 2009; Kushnir et al., 2006; Gong et al., 2012; Wang et al., 2011). As a first stage in their semi-supervised learning method, Goldberg, Zhu, Singh, Xu, and Nowak (2009) design a spectral clustering algorithm. The method starts by subsampling the data points, obtaining ‘centers’ in the following way. Draw  $\mathbf{y}_1$  at random from the data and remove its  $\ell$ -nearest neighbors from the data. Then repeat with the remaining data, obtaining centers  $\mathbf{y}_1, \mathbf{y}_2, \dots$ . Let  $\mathbf{C}_i$  denote the sample covariance in the neighborhood of  $\mathbf{y}_i$  made of its  $\ell$ -nearest neighbors. An  $m$ -nearest-neighbor graph is then defined on the centers in terms of the Mahalanobis distances. Explicitly, the centers  $\mathbf{y}_i$  and  $\mathbf{y}_j$  are connected in the graph if  $\mathbf{y}_j$  is among the  $m$  nearest neighbors of  $\mathbf{y}_i$  in Mahalanobis distance

$$\|\mathbf{C}_i^{-1/2}(\mathbf{y}_i - \mathbf{y}_j)\|, \tag{1}$$

or vice-versa. The parameters  $\ell$  and  $m$  are both chosen of order  $\log n$ . An existing edge between  $\mathbf{y}_i$  and  $\mathbf{y}_j$  is then weighted by  $\exp(-H_{ij}^2/\eta^2)$ , where  $H_{ij}$  denotes the Hellinger distance between the probability distributions  $\mathcal{N}(\mathbf{0}, \mathbf{C}_i)$  and  $\mathcal{N}(\mathbf{0}, \mathbf{C}_j)$ . The spectral graph partitioning algorithm of Ng, Jordan, and Weiss (2002)—detailed in Algorithm 1—is then applied to the resulting affinity matrix, with some form of constrained K-means. We note that Goldberg et al. (2009) evaluate their method in the context of semi-supervised learning where the clustering routine is only required to return subclusters of actual clusters. In particular, the data points other than the centers are discarded. Note also that their evaluation is empirical.

---

**Algorithm 1** Spectral Graph Partitioning (Ng, Jordan, and Weiss, 2002)

---

**Input:**

Affinity matrix  $\mathbf{W} = (W_{ij})$ , size of the partition  $K$

**Steps:**

- 1: Compute  $\mathbf{Z} = (Z_{ij})$  according to  $Z_{ij} = W_{ij}/\sqrt{\Delta_i\Delta_j}$ , with  $\Delta_i := \sum_{j=1}^n W_{ij}$ .
  - 2: Extract the top  $K$  eigenvectors of  $\mathbf{Z}$ .
  - 3: Renormalize each row of the resulting  $n \times K$  matrix.
  - 4: Apply  $K$ -means to the row vectors.
- 

The algorithm proposed by Kushnir, Galun, and Brandt (2006) is multiscale and works by coarsening the neighborhood graph and computing sampling density and geometric information inferred along the way such as obtained via PCA in local neighborhoods. This bottom-up flow is then followed by a top-down pass, and the two are iterated a few times. The algorithm is too complex to be described in detail here, and probably too complex to be analyzed mathematically. The clustering methods of Goldberg et al. (2009) and ours can be seen as simpler variants that only go bottom up and coarsen the graph only once.

In the last stages of writing this paper, we learned of the works of Wang, Jiang, Wu, and Zhou (2011) and Gong, Zhao, and Medioni (2012), who propose algorithms very similar to our Algorithm 3 detailed below. Note that these publications do not provide any theoretical guarantees for their methods, which is one of our main contributions here.

## 2.2 Our Algorithms

We now describe our method and propose several variants. Our setting is standard: we observe data points  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^D$  that we assume were sampled in the vicinity of  $K$  smooth surfaces embedded in  $\mathbb{R}^D$ . The setting is formalized later in Section 3.1.

### 2.2.1 CONNECTED COMPONENT EXTRACTION: COMPARING LOCAL COVARIANCES

We start with our simplest variant, which is also the most natural. The method depends on a neighborhood radius  $r > 0$ , a spatial scale parameter  $\varepsilon > 0$  and a covariance (relative) scale  $\eta > 0$ . For a vector  $\mathbf{x}$ ,  $\|\mathbf{x}\|$  denotes its Euclidean norm, and for a (square) matrix  $\mathbf{A}$ ,  $\|\mathbf{A}\|$  denotes its spectral norm. For  $n \in \mathbb{N}$ , we denote by  $[n]$  the set  $\{1, \dots, n\}$ . Given a data set  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , for any point  $\mathbf{x} \in \mathbb{R}^D$  and  $r > 0$ , define the neighborhood

$$N_r(\mathbf{x}) = \{\mathbf{x}_j : \|\mathbf{x} - \mathbf{x}_j\| \leq r\}.$$

---

#### Algorithm 2 Connected Component Extraction: Comparing Covariances

---

**Input:**

Data points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ ; neighborhood radius  $r > 0$ ; spatial scale  $\varepsilon > 0$ , covariance scale  $\eta > 0$ .

**Steps:**

- 1: For each  $i \in [n]$ , compute the sample covariance matrix  $\mathbf{C}_i$  of  $N_r(\mathbf{x}_i)$ .
- 2: Remove  $\mathbf{x}_i$  when there is  $\mathbf{x}_j$  such that  $\|\mathbf{x}_j - \mathbf{x}_i\| \leq r$  and  $\|\mathbf{C}_j - \mathbf{C}_i\| > \eta r^2$ .
- 3: Compute the following affinities between data points:

$$W_{ij} = \mathbb{I}_{\{\|\mathbf{x}_i - \mathbf{x}_j\| \leq \varepsilon\}} \cdot \mathbb{I}_{\{\|\mathbf{C}_i - \mathbf{C}_j\| \leq \eta r^2\}}. \quad (2)$$

- 4: Extract the connected components of the resulting graph.
  - 5: Each point removed in Step 2 is grouped with the closest point that survived Step 2.
- 

In words, the algorithm first computes local covariances (Step 1). It removes points that are believed to be very near an intersection (Step 2)—we elaborate on this below. With the remaining points, it creates an unweighted graph: the nodes of this graph are the data points and edges are formed between two nodes if both the distance between these nodes and the distance between the local covariance structures at these nodes are sufficiently small (Step 3). The connected components of the resulting graph are extracted (Step 4) and the points that survived Step 2 are labeled accordingly. Each point removed in Step 2 receives the label of its closest labeled point (Step 5).

In principle, the neighborhood size  $r$  is chosen just large enough that performing PCA in each neighborhood yields a reliable estimate of the local covariance structure. For this, the number of points inside the neighborhood needs to be large enough, which depends on the sample size  $n$ , the sampling density, intrinsic dimension of the surfaces and their surface area (Hausdorff measure), how far the points are from the surfaces (i.e., noise level), and the regularity of the surfaces. The spatial scale parameter  $\varepsilon$  depends on the sampling density and  $r$ . It is meant to be larger than  $r$  and needs to be large enough that a point has plenty

of points within distance  $\varepsilon$ , including some across an intersection, so each cluster is strongly connected. At the same time,  $\varepsilon$  needs to be small enough that a local linear approximation to the surfaces is a relevant feature of proximity. Its choice is rather similar to the choice of the scale parameter in standard spectral clustering (Ng et al., 2002; Zelnik-Manor and Perona, 2005). The covariance scale  $\eta$  needs to be large enough that centers from the same cluster and within distance  $\varepsilon$  of each other have local covariance matrices within distance  $\eta r^2$ , but small enough that points from different clusters near their intersection have local covariance matrices separated by a distance substantially larger than  $\eta r^2$ . This depends on the curvature of the surfaces and the incidence angle at the intersection of two (or more) surfaces. Note that a typical covariance matrix over a ball of radius  $r$  has norm of order  $r^2$ , which justifies using our choice of parametrization. In the mathematical framework we introduce later on, these parameters can be chosen automatically as done in (Arias-Castro et al., 2011), at least when the points are sampled exactly on the surfaces. We will not elaborate on that since in practice this does not inform our choice of parameters.

The rationale behind Step 2 is as follows. As we just discussed, the parameters need to be tuned so that points from the same cluster and within distance  $\varepsilon$  have local covariance matrices within distance  $\eta r^2$ . Strictly speaking, this is true of points away from the boundary of the underlying surface. Hence, although  $\mathbf{x}_i$  and  $\mathbf{x}_j$  in Step 2 are meant to be from different clusters, they could be from the same surface near its boundary. In the former situation, since they are near each other, in our model this will imply that they are close to an intersection. Therefore, roughly speaking, Step 2 removes points near an intersection, but also points near the boundaries of the underlying surfaces. The reason that we require  $\mathbf{x}_i$  and  $\mathbf{x}_j$  to be within distance  $r$ , as opposed to  $\varepsilon$ , is because otherwise removing the points in Step 2 would create a “gap” of  $\varepsilon$  near an intersection which then cannot be bridged in Steps 3-4. (Alternatively, one could replace  $r$  with  $\xi \in (r, \varepsilon)$  in Step 2, but this would add this additional parameter  $\xi$  to the algorithm.) This step is in fact crucial as the local covariance varies smoothly along the intersection of two smooth surfaces.

Although this method works in simple situations like that of two intersecting segments (Figure 1), it is not meant to be practical. Indeed, extracting connected components is known to be sensitive to spurious points and therefore unstable. Furthermore, we found that comparing local covariance matrices as in affinity (2) tends to be less stable than comparing local projections as in affinity (3), which brings us to our next variant.

### 2.2.2 CONNECTED COMPONENT EXTRACTION: COMPARING LOCAL PROJECTIONS

We present another variant also based on extracting the connected components of a neighborhood graph that compares orthogonal projections onto the largest principal directions. See Algorithm 3.

We note that the local intrinsic dimension is determined by thresholding the eigenvalues of the local covariance matrix, keeping the directions with eigenvalues within some range of the largest eigenvalue. The same strategy is used by Kushnir et al. (2006), but with a different threshold. The method is a hard version of what we implemented, which we describe in Algorithm 4.

We note that neither algorithm includes an intersection-removal step as Step 2 in Algorithm 2. The reason is that the algorithms work without such a step. Indeed, we show in

---

**Algorithm 3** Connected Component Extraction: Comparing Projections

---

**Input:**

Data points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ ; neighborhood radius  $r > 0$ , spatial scale  $\varepsilon > 0$ , projection scale  $\eta > 0$ .

**Steps:**

- 1:** For each  $i \in [n]$ , compute the sample covariance matrix  $\mathbf{C}_i$  of  $N_r(\mathbf{x}_i)$ .
- 2:** Compute the projection  $\mathbf{Q}_i$  onto the eigenvectors of  $\mathbf{C}_i$  with corresponding eigenvalue exceeding  $\sqrt{\eta} \|\mathbf{C}_i\|$ .
- 3:** Compute the following affinities between data points:

$$W_{ij} = \mathbb{I}_{\{\|\mathbf{x}_i - \mathbf{x}_j\| \leq \varepsilon\}} \cdot \mathbb{I}_{\{\|\mathbf{Q}_i - \mathbf{Q}_j\| \leq \eta\}}. \quad (3)$$

- 4:** Extract the connected components of the resulting graph.
- 

Theorem 1 that Algorithm 3 is able to separate the clusters—the only drawback is that it may possibly treat the intersection of two surfaces as a cluster. And we show via numerical experiments in Section 4 that Algorithm 4 performs well in a number of situations.

## 2.2.3 COVARIANCES OR PROJECTIONS?

In our numerical experiments, we tried working both directly with covariance matrices as in (2) and with projections as in (3). Note that in our experiments we used spectral graph partitioning with soft versions of these affinities, as described in Section 2.2.4. We found working with projections to be more reliable. The problem comes, in part, from boundaries. When a surface has a boundary, local covariances over neighborhoods that overlap with the boundary are quite different from local covariances over nearby neighborhoods that do not touch the boundary. Consider the example of two segments,  $S_1$  and  $S_2$ , intersecting at an angle of  $\theta \in (0, \pi/2)$  at their middle point, specifically

$$S_1 = [-1, 1] \times \{0\}, \quad S_2 = \{(x, x \tan \theta) : x \in [-\cos \theta, \cos \theta]\}.$$

Assume there is no noise and that the sampling is uniform. Assume  $r \in (0, \frac{1}{2} \sin \theta)$  so that the disc centered at  $\mathbf{x}_1 := (1/2, 0)$  does not intersect  $S_2$ , and the disc centered at  $\mathbf{x}_2 := (\frac{1}{2} \cos \theta, \frac{1}{2} \sin \theta)$  does not intersect  $S_1$ . Let  $\mathbf{x}_0 = (1, 0)$ . For  $\mathbf{x} \in S_1 \cup S_2$ , let  $\mathbf{C}_\mathbf{x}$  denote the local covariance at  $\mathbf{x}$  over a ball of radius  $r < \frac{1}{2} \sin \theta$ , so that the neighborhoods of  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2$  are “pure”. The situation is described in Figure 2.

Simple calculations yield

$$\mathbf{C}_{\mathbf{x}_0} = \frac{r^2}{12} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbf{C}_{\mathbf{x}_1} = \frac{r^2}{3} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbf{C}_{\mathbf{x}_2} = \frac{r^2}{3} \begin{pmatrix} \cos^2 \theta & \sin(\theta) \cos(\theta) \\ \sin(\theta) \cos(\theta) & \sin^2 \theta \end{pmatrix},$$

so that

$$\|\mathbf{C}_{\mathbf{x}_0} - \mathbf{C}_{\mathbf{x}_1}\| = \frac{r^2}{4}, \quad \|\mathbf{C}_{\mathbf{x}_1} - \mathbf{C}_{\mathbf{x}_2}\| = \frac{r^2}{3} \sin \theta.$$

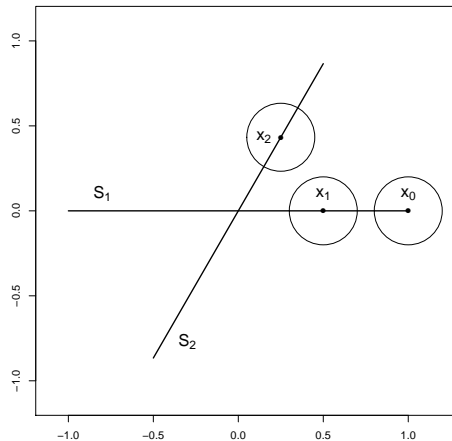


Figure 2: Two segments intersecting. The local covariances (within the disc neighborhoods drawn) at  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are closer than the local covariances at  $\mathbf{x}_1$  and  $\mathbf{x}_0$ , even though  $\mathbf{x}_0$  and  $\mathbf{x}_1$  are on the same segment.

Therefore, when  $\sin \theta \leq \frac{3}{4}$  (roughly,  $\theta \leq 48^\circ$ ), the local covariances at  $\mathbf{x}_0, \mathbf{x}_1 \in S_1$  are farther (in operator norm) than those at  $\mathbf{x}_1 \in S_1$  and  $\mathbf{x}_2 \in S_2$ . As for projections, however,

$$\mathbf{Q}_{\mathbf{x}_0} = \mathbf{Q}_{\mathbf{x}_1} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbf{Q}_{\mathbf{x}_2} = \begin{pmatrix} \cos^2 \theta & \sin(\theta) \cos(\theta) \\ \sin(\theta) \cos(\theta) & \sin^2 \theta \end{pmatrix},$$

so that

$$\|\mathbf{Q}_{\mathbf{x}_0} - \mathbf{Q}_{\mathbf{x}_1}\| = 0, \quad \|\mathbf{Q}_{\mathbf{x}_1} - \mathbf{Q}_{\mathbf{x}_2}\| = \sqrt{2} \sin \theta.$$

While in theory points within distance  $r$  from the boundary account for a small portion of the sample, in practice this is not the case, at least not with the sample sizes that we are able to rapidly process. In fact, we find that spectral graph partitioning is challenged by having points near the boundary that are far (in affinity) from nearby points from the same cluster. This may explain why the (soft version of) affinity (3) yields better results than the (soft version of) affinity (2) in our experiments.

#### 2.2.4 SPECTRAL CLUSTERING BASED ON LOCAL PCA

The following variant is more robust in practice and is the algorithm we actually implemented. The method assumes that the surfaces are of same dimension  $d$  and that there are  $K$  of them, with both parameters  $K$  and  $d$  known.

We note that  $\mathbf{y}_1, \dots, \mathbf{y}_{n_0}$  forms an  $r$ -packing of the data. The underlying rationale for this coarsening is justified in (Goldberg et al., 2009) by the fact that the covariance matrices, and also the top principal directions, change smoothly with the location of the neighborhood, so that without subsampling these characteristics would not help detect the abrupt event of an intersection. The affinity (4) is of course a soft version of (3).



---

**Algorithm 4** Spectral Clustering Based on Local PCA

---

**Input:**

Data points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ ; neighborhood radius  $r > 0$ ; spatial scale  $\varepsilon > 0$ , projection scale  $\eta > 0$ ; intrinsic dimension  $d$ ; number of clusters  $K$ .

**Steps:**

**0:** Pick one point  $\mathbf{y}_1$  at random from the data. Pick another point  $\mathbf{y}_2$  among the data points not included in  $N_r(\mathbf{y}_1)$ , and repeat the process, selecting centers  $\mathbf{y}_1, \dots, \mathbf{y}_{n_0}$ .

**1:** For each  $i = 1, \dots, n_0$ , compute the sample covariance matrix  $\mathbf{C}_i$  of  $N_r(\mathbf{y}_i)$ . Let  $\mathbf{Q}_i$  denote the orthogonal projection onto the space spanned by the top  $d$  eigenvectors of  $\mathbf{C}_i$ .

**2:** Compute the following affinities between center pairs:

$$W_{ij} = \exp\left(-\frac{\|\mathbf{y}_i - \mathbf{y}_j\|^2}{\varepsilon^2}\right) \cdot \exp\left(-\frac{\|\mathbf{Q}_i - \mathbf{Q}_j\|^2}{\eta^2}\right). \quad (4)$$

**3:** Apply spectral graph partitioning (Algorithm 1) to  $\mathbf{W}$ .

**4:** The data points are clustered according to the closest center in Euclidean distance.

---

## 2.2.5 COMPARISON WITH CLOSELY RELATED METHODS

We highlight some differences with the other proposals in the literature. We first compare our approach to that of Goldberg et al. (2009), which was our main inspiration.

- *Neighborhoods.* Comparing with Goldberg et al. (2009), we define neighborhoods over  $r$ -balls instead of  $\ell$ -nearest neighbors, and connect points over  $\varepsilon$ -balls instead of  $m$ -nearest neighbors. This choice is for convenience, as these ways are in fact essentially equivalent when the sampling density is fairly uniform. This is elaborated at length in (Maier et al., 2009; Brito et al., 1997; Arias-Castro, 2011).
- *Mahalanobis distances.* Goldberg et al. (2009) use Mahalanobis distances (1) between centers. In our version, we could for example replace the Euclidean distance  $\|\mathbf{x}_i - \mathbf{x}_j\|$  in the affinity (2) with the average Mahalanobis distance

$$\|\mathbf{C}_i^{-1/2}(\mathbf{x}_i - \mathbf{x}_j)\| + \|\mathbf{C}_j^{-1/2}(\mathbf{x}_j - \mathbf{x}_i)\|.$$

We actually tried this and found that the algorithm was less stable, particularly under low noise. Introducing a regularization in this distance—which requires the introduction of another parameter—solves this problem, at least partially.

That said, using Mahalanobis distances makes the procedure less sensitive to the choice of  $\varepsilon$ , in that neighborhoods may include points from different clusters. Think of two parallel line segments separated by a distance of  $\delta$ , and assume there is no noise, so the points are sampled exactly from these segments. Assuming an infinite sample size, the local covariance is the same everywhere so that points within distance  $\varepsilon$  are connected by the affinity (2). Hence, Algorithm 2 requires that  $\varepsilon < \delta$ . In terms of Mahalanobis distances, points on different segments are infinitely separated, so a

version based on these distances would work with any  $\varepsilon > 0$ . In the case of curved surfaces and/or noise, the situation is similar, though not as evident. Even then, the gain in performance is not obvious since we only require that  $\varepsilon$  be slightly larger in order of magnitude than  $r$ .

- *Hellinger distances.* As we mentioned earlier, Goldberg et al. (2009) use Hellinger distances of the probability distributions  $\mathcal{N}(\mathbf{0}, \mathbf{C}_i)$  and  $\mathcal{N}(\mathbf{0}, \mathbf{C}_j)$  to compare covariance matrices, specifically

$$\left( 1 - 2^{D/2} \frac{\det(\mathbf{C}_i \mathbf{C}_j)^{1/4}}{\det(\mathbf{C}_i + \mathbf{C}_j)^{1/2}} \right)^{1/2}, \quad (5)$$

if  $\mathbf{C}_i$  and  $\mathbf{C}_j$  are full-rank. While using these distances or the Frobenius distances makes little difference in practice, we find it easier to work with the latter when it comes to proving theoretical guarantees. Moreover, it seems more natural to assume a uniform sampling distribution in each neighborhood rather than a normal distribution, so that using the more sophisticated similarity (5) does not seem justified.

- *K-means.* We use K-means++ for a good initialization. We found that the more sophisticated size-constrained K-means (Bradley et al., 2000) used in (Goldberg et al., 2009) did not improve the clustering results.

As we mentioned above, our work was developed in parallel to that of Wang et al. (2011) and Gong et al. (2012). We highlight some differences. First, there is no subsampling, but rather, the local tangent space is estimated at each data point  $\mathbf{x}_i$ . Wang et al. (2011) fit a mixture of  $d$ -dimensional affine subspaces to the data using MPPCA (Tipping and Bishop, 1999), which is then used to estimate the tangent subspaces at each data point. Gong et al. (2012) develop some sort of robust local PCA. While Wang et al. (2011) assume all surfaces are of same dimension known to the user, Gong et al. (2012) estimate that locally by looking at the largest gap in the spectrum of estimated local covariance matrix. This is similar in spirit to what is done in Step 2 of Algorithm 3, but we did not include this step in Algorithm 4 because we did not find it reliable in practice. We also tried estimating the local dimensionality using the method of Little et al. (2009), but this failed in the most complex cases.

Wang et al. (2011) use a nearest-neighbor graph and their affinity is defined as

$$W_{ij} = \Delta_{ij} \cdot \left( \prod_{s=1}^d \cos \theta_s(i, j) \right)^\alpha,$$

where  $\Delta_{ij} = 1$  if  $\mathbf{x}_i$  is among the  $\ell$ -nearest neighbors of  $\mathbf{x}_j$ , or vice versa, while  $\Delta_{ij} = 0$  otherwise;  $\theta_1(i, j) \geq \dots \geq \theta_d(i, j)$  are the principal (a.k.a., canonical) angles (Stewart and Sun, 1990) between the estimated tangent subspaces at  $\mathbf{x}_i$  and  $\mathbf{x}_j$ .  $\ell$  and  $\alpha$  are parameters of the method. Gong et al. (2012) define an affinity that incorporates the self-tuning method of Zelnik-Manor and Perona (2005); in our notation, their affinity is

$$\exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\varepsilon_i \varepsilon_j}\right) \cdot \exp\left(-\frac{(\sin^{-1}(\|\mathbf{Q}_i - \mathbf{Q}_j\|))^2}{\eta^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2 / (\varepsilon_i \varepsilon_j)}\right).$$

where  $\varepsilon_i$  is the distance from  $\mathbf{x}_i$  to its  $\ell$ -nearest neighbor.  $\ell$  is a parameter.

Although we do not analyze their respective ways of estimating the tangent subspaces, our analysis provides essential insights into their methods, and for that matter, any other method built on spectral clustering based on tangent subspace comparisons.

### 3. Mathematical Analysis

While the analysis of Algorithm 4 seems within reach, there are some complications due to the fact that points near the intersection may form a cluster of their own—we were not able to discard this possibility. Instead, we study the simpler variants described in Algorithm 2 and Algorithm 3. Even then, the arguments are rather complex and interestingly involved. The theoretical guarantees that we obtain for these variants are stated in Theorem 1 and proved in Section 6. We comment on the analysis of Algorithm 4 right after that. We note that there are very few theoretical results on resolving intersecting manifolds—in fact, we are only aware of (Arias-Castro et al., 2011) (under severe restrictions on the dimension of the intersection). Some such results have been established for a number of methods for subspace clustering (affine surfaces), for example, in (Chen and Lerman, 2009b; Soltanolkotabi and Candès, 2012; Soltanolkotabi et al., 2014; Wang and Xu, 2013; Heckel and Bölcskei, 2013; Tsakiris and Vidal, 2015; Ma et al., 2008).

The generative model we assume is a natural mathematical framework for multi-manifold learning where points are sampled in the vicinity of smooth surfaces embedded in Euclidean space. For concreteness and ease of exposition, we focus on the situation where two surfaces (i.e.,  $K = 2$ ) of same dimension  $1 \leq d \leq D$  intersect. This special situation already contains all the geometric intricacies of separating intersecting clusters. On the one hand, clusters of different intrinsic dimension may be separated with an accurate estimation of the local intrinsic dimension without further geometry involved (Haro et al., 2007). On the other hand, more complex intersections (3-way and higher) complicate the situation without offering truly new challenges. For simplicity of exposition, we assume that the surfaces are submanifolds without boundary, though it will be clear from the analysis (and the experiments) that the method can handle surfaces with (smooth) boundaries that may self-intersect. We discuss other possible extensions in Section 5.

Within that framework, we show that Algorithm 2 and Algorithm 3 are able to identify the clusters accurately except for points near the intersection. Specifically, with high probability with respect to the sampling distribution, Algorithm 2 divides the data points into two groups such that, except for points within distance  $C\varepsilon$  of the intersection, all points from the first cluster are in one group and all points from the second cluster are in the other group. The constant  $C$  depends on the surfaces, including their curvatures, separation between them and intersection angle. The situation for Algorithm 3 is more complex, as it may return more than two clusters, but most of the two clusters (again, away from the intersection) are in separate connected components.

#### 3.1 Generative Model

Each surface we consider is a connected,  $C^2$  and compact submanifold without boundary and of dimension  $d$  embedded in  $\mathbb{R}^D$ . Any such surface has a positive reach, which is what we use to quantify smoothness. The notion of reach was introduced by Federer (1959).

Intuitively, a surface has reach exceeding  $r$  if and only if one can roll a ball of radius  $r$  on the surface without obstruction (Walther, 1997; Cuevas et al., 2012). Formally, for  $\mathbf{x} \in \mathbb{R}^D$  and  $S \subset \mathbb{R}^D$ , let

$$\text{dist}(\mathbf{x}, S) = \inf_{\mathbf{s} \in S} \|\mathbf{x} - \mathbf{s}\|,$$

and

$$B(S, r) = \{\mathbf{x} : \text{dist}(\mathbf{x}, S) < r\},$$

which is often called the  $r$ -tubular neighborhood (or  $r$ -neighborhood) of  $S$ . The reach of  $S$  is the supremum over  $r > 0$  such that, for each  $\mathbf{x} \in B(S, r)$ , there is a unique point in  $S$  nearest  $\mathbf{x}$ . It is well-known that, for  $C^2$  submanifolds, the reach bounds the radius of curvature from below (Federer, 1959, Lem. 4.17). For submanifolds without boundaries, the reach coincides with the condition number introduced in (Niyogi et al., 2008).

When two surfaces  $S_1$  and  $S_2$  intersect, meaning  $S_1 \cap S_2 \neq \emptyset$ , we define their incidence angle as

$$\theta(S_1, S_2) := \inf \{\theta_{\min}(T_{S_1}(\mathbf{s}), T_{S_2}(\mathbf{s})) : \mathbf{s} \in S_1 \cap S_2\}, \quad (6)$$

where  $T_S(\mathbf{s})$  denote the tangent subspace of submanifold  $S$  at point  $\mathbf{s} \in S$ , and  $\theta_{\min}(T_1, T_2)$  is the smallest *nonzero* principal (a.k.a., canonical) angle between subspaces  $T_1$  and  $T_2$  (Stewart and Sun, 1990).

The clusters are generated as follows. Each data point  $\mathbf{x}_i$  is drawn according to

$$\mathbf{x}_i = \mathbf{s}_i + \mathbf{z}_i, \quad (7)$$

where  $\mathbf{s}_i$  is drawn from the uniform distribution over  $S_1 \cup S_2$  and  $\mathbf{z}_i$  is an additive noise term satisfying  $\|\mathbf{z}_i\| \leq \tau$ —thus  $\tau$  represents the noise or jitter level. When  $\tau = 0$  the points are sampled exactly on the surfaces. We assume the points are sampled independently of each other. We let

$$I_k = \{i : \mathbf{s}_i \in S_k\},$$

and the goal is to recover the groups  $I_1$  and  $I_2$ , up to some errors.

### 3.2 Performance Guarantees

We state some performance guarantees for Algorithm 2 and Algorithm 3.

**Theorem 1.** *Consider two connected, compact, twice continuously differentiable submanifolds without boundary, of same dimension  $d$ , intersecting at a strictly positive angle, with the intersection set having strictly positive reach. Assume the parameters are set so that*

$$\tau \leq r\eta/C, \quad r \leq \varepsilon/C, \quad \varepsilon \leq \eta/C, \quad \eta \leq 1/C, \quad (8)$$

for a large-enough constant  $C \geq 1$  that depends on the configuration. Then with probability at least  $1 - Cn \exp[-nr^d\eta^2/C]$ :

- Algorithm 2 returns exactly two groups such that two points from different clusters are not grouped together unless one of them is within distance  $Cr$  from the intersection.
- Algorithm 3 returns at least two groups, and such that two points from different clusters are not grouped together unless one of them is within distance  $Cr$  from the intersection.

Thus, as long as, (8) is satisfied, the algorithms have the above properties with high probability when  $r^d \eta^2 \geq C' \log n/n$ , where  $C' > C$  is fixed. In particular, we may choose  $\eta \asymp 1$  and  $\varepsilon \asymp r \asymp \tau \vee (\log(n)/n)^{1/d}$ , which lightens the computational burden.

We note that, while the constant  $C > 0$  does not depend on the sample size  $n$ , it depends in somewhat complicated ways on characteristics of the surfaces and their position relative to each other, such as their reach and intersection angle, but also aspects that are harder to quantify, like their separation away from their intersection. We note, however, that it behaves as expected:  $C$  is indeed decreasing in the reach and the intersection angle, and increasing in the intrinsic dimensions of the surfaces, for example.

The algorithms may make clustering mistakes within distance  $Cr$  of the intersection, where  $Cr \asymp \tau \vee (\log(n)/n)^{1/d}$  with the choice of parameters just described. Whether this is optimal in the nonparametric setting that we consider—for example, in a minimax sense—we do not know.

We now comment on the challenge of proving a similar result for Algorithm 4. This algorithm relies on knowledge of the intrinsic dimension of the surfaces  $d$  and the number of clusters (here  $K = 2$ ), but these may be estimated as in (Arias-Castro et al., 2011), at least in theory, so we assume these parameters are known. The subsampling done in Step 0 does not pose any problem whatsoever, since the centers are well-spread when the points themselves are. The difficulty resides in the application of the spectral graph partitioning, Algorithm 1. If we were to include the intersection-removal step (Step 2 of Algorithm 2) before applying spectral graph partitioning, then a simple adaptation of arguments in (Arias-Castro, 2011) would suffice. The real difficulty, and potential pitfall of the method in this framework (without the intersection-removal step), is that the points near the intersection may form their own cluster. For example, in the simplest case of two affine surfaces intersecting at a positive angle and no sampling noise, the projection matrix at a point near the intersection—meaning a point whose  $r$ -ball contains a substantial piece of both surfaces—would be the projection matrix onto  $S_1 + S_2$  seen as a linear subspace. We were not able to discard this possibility, although we do not observe this happening in practice. A possible remedy is to constrain the K-means part to only return large-enough clusters. However, a proper analysis of this would require a substantial amount of additional work and we did not engage seriously in this pursuit.

## 4. Numerical Experiments

### 4.1 Some Illustrative Examples

We started by applying our method<sup>2</sup> on a few artificial examples to illustrate the theory. As we argued earlier, the methods of Wang et al. (2011) and Gong et al. (2012) are quite similar to ours, and we encourage the reader to also look at the numerical experiments they performed. Our numerical experiments should be regarded as a proof of concept, only here to show that our method can be implemented and works on some stylized examples.

In all experiments, the number of clusters  $K$  and the dimension of the manifolds  $d$  are assumed known. We choose the spatial scale  $\varepsilon$  and the projection scale  $\eta$  automatically as

---

2. The code is available online at <https://math.cos.ucf.edu/tengz>.

follows: we let

$$\varepsilon = \max_{1 \leq i \leq n_0} \min_{j \neq i} \|\mathbf{y}_i - \mathbf{y}_j\|, \tag{9}$$

and

$$\eta = \operatorname{median}_{(i,j): \|\mathbf{y}_i - \mathbf{y}_j\| < \varepsilon} \|\mathbf{Q}_i - \mathbf{Q}_j\|.$$

Here, we implicitly assume that the union of all the underlying surfaces forms a connected set. In that case, the idea behind choosing  $\varepsilon$  as in (9) is that we want the  $\varepsilon$ -graph on the centers  $\mathbf{y}_1, \dots, \mathbf{y}_n$  to be connected. Then  $\eta$  is chosen so that a center  $\mathbf{y}_i$  remains connected in the  $(\varepsilon, \eta)$ -graph to most of its neighbors in the  $\varepsilon$ -graph.

The neighborhood radius  $r$  is chosen by hand for each situation. Although we do not know how to choose  $r$  automatically, there are some general ad hoc guidelines. When  $r$  is too large, the local linear approximation to the underlying surfaces may not hold in neighborhoods of radius  $r$ , resulting in local PCA becoming inappropriate. When  $r$  is too small, there might not be enough points in a neighborhood of radius  $r$  to accurately estimate the local tangent subspace to a given surface at that location, resulting in local PCA becoming inaccurate. From a computational point of view, the smaller  $r$ , the larger the number of neighborhoods and the heavier the computations, particularly at the level of spectral graph partitioning. In our numerical experiments, we find that our algorithm is more sensitive to the choice of  $r$  when the clustering problem is more difficult. We note that automatic choice of tuning parameters remains a challenge in clustering, and machine learning at large, especially when no labels are available whatsoever. See (Zelnik-Manor and Perona, 2005; Zhang et al., 2012; Little et al., 2009; Kaslovsky and Meyer, 2011).

Since the algorithm is randomized (see Step 0 in Algorithm 4) we repeat each simulation 100 times and report the median misclustering rate and number of times where the misclustering rate is smaller than 5%, 10%, and 15%.

We first run Algorithm 4 on several artificial data sets, which are demonstrated in the LHS of Figures 3 and 4. Table 1 reports the local radius  $r$  used for each data set ( $R$  is the global radius of each data set), and the statistics for misclustering rates. Typical clustering results are demonstrated in the RHS of Figures 3 and 4. It is evident that Algorithm 4 performs well in these simulations.

data set	$r$	median misclustering rate	5%	10%	15%
Three curves	0.02 (0.034 <i>R</i> )	4.16%	76	89	89
Self-intersecting curves	0.1 (0.017 <i>R</i> )	1.16%	85	85	86
Two spheres	0.2 (0.059 <i>R</i> )	3.98%	100	100	100
Möbius strips	0.1 (0.028 <i>R</i> )	2.22%	85	86	88
Monkey saddle	0.1 (0.069 <i>R</i> )	9.73%	0	67	97
Paraboloids	0.07 (0.048 <i>R</i> )	10.42%	0	12	91

Table 1: Choices for  $r$  and misclustering statistics for the artificial data sets demonstrated in Figures 3 and 4. The statistics are based on 100 repeats and include the median misclustering rate and number of repeats where the misclustering rate is smaller than 5%, 10% and 15%.

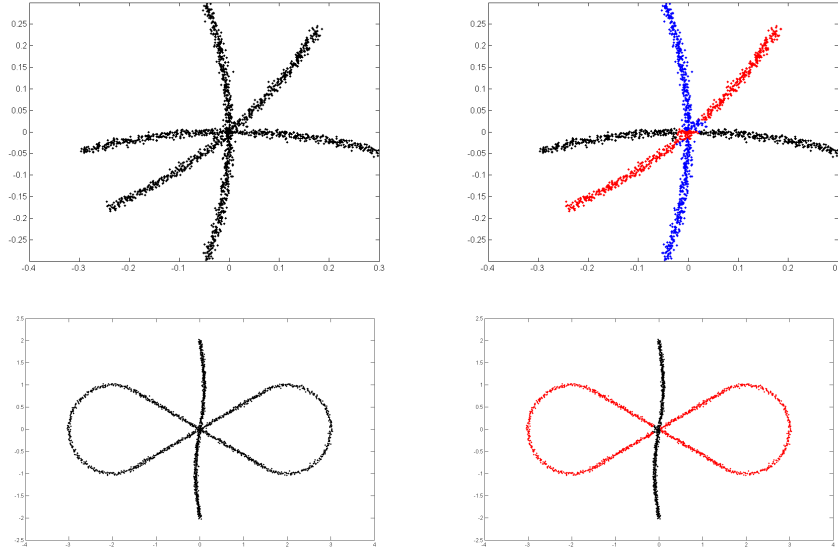


Figure 3: Performance of Algorithm 4 on data sets “Three curves” and “Self-intersecting curves”. Left column is the input data sets, and right column demonstrates the typical clustering.

In another simulation, we show the dependence of the success of our algorithm on the intersecting angle between curves in Table 2 and Figure 5. Here, we fix two curves intersecting at a point, and gradually decrease the intersection angle by rotating one of them while holding the other one fixed. The angles are  $\pi/2$ ,  $\pi/4$ ,  $\pi/6$  and  $\pi/8$ . From the table we can see that our algorithm performs well when the angle is  $\pi/4$ , but the performance deteriorates as the angle becomes smaller, and the algorithm almost always fails when the angle is  $\pi/8$ .

Intersecting angle	$r$	median misclustering rate	5%	10%	15%
$\pi/2$	0.02 (0.034 <i>R</i> )	2.08%	98	98	98
$\pi/4$	0.02 (0.034 <i>R</i> )	3.33%	92	94	94
$\pi/6$	0.02 (0.034 <i>R</i> )	5.53%	32	59	59
$\pi/8$	0.02 (0.033 <i>R</i> )	27.87%	0	2	2

Table 2: Choices for  $r$  and misclustering statistics for the instances of two intersecting curves demonstrated in Figure 5. The statistics are based on 100 repeats and include the median misclustering rate and number of repeats where the misclustering rate is smaller than 5%, 10% and 15%.

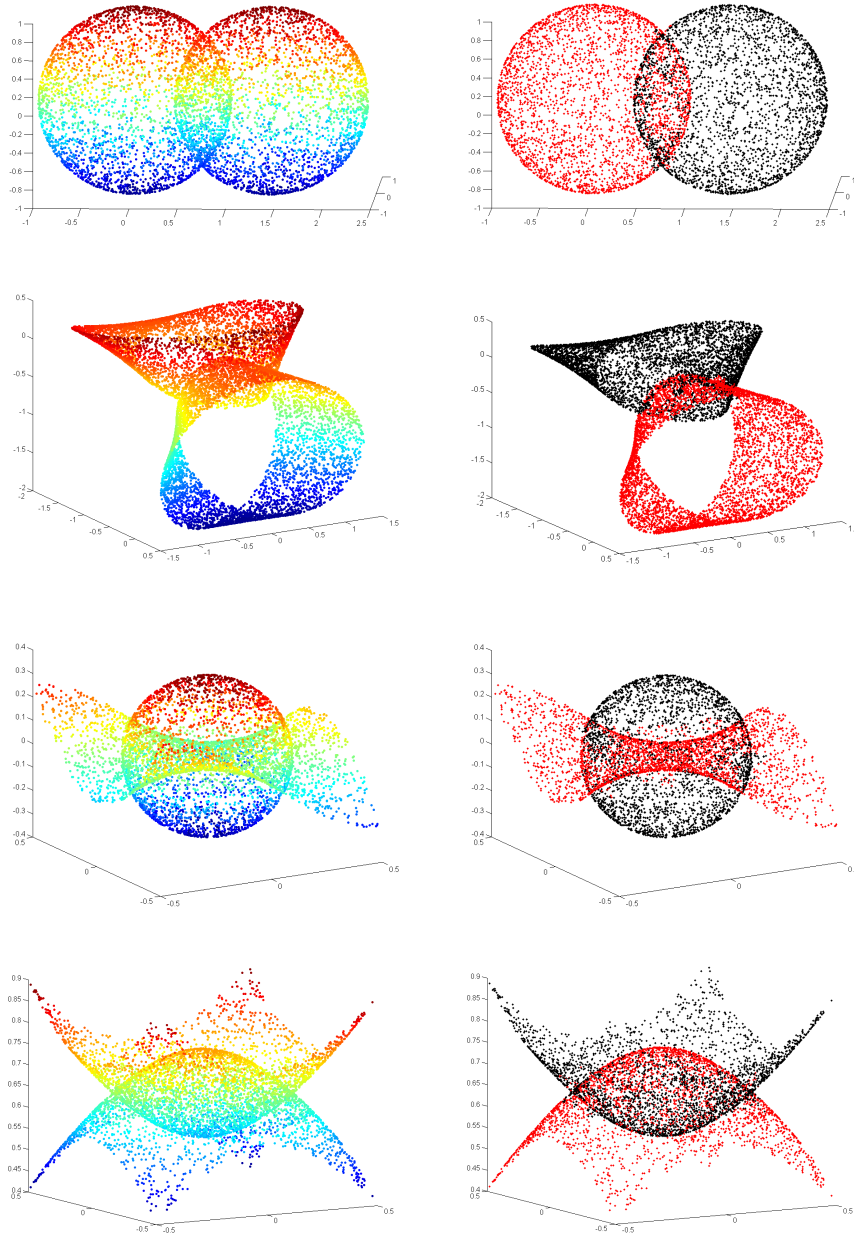


Figure 4: Performance of Algorithm 4 on data sets “Two spheres”, “Möbius strips”, “Monkey saddle” and “Paraboloids”. Left column is the input data sets, and right column demonstrates the typical clustering.



# SPECTRAL CLUSTERING BASED ON LOCAL PCA

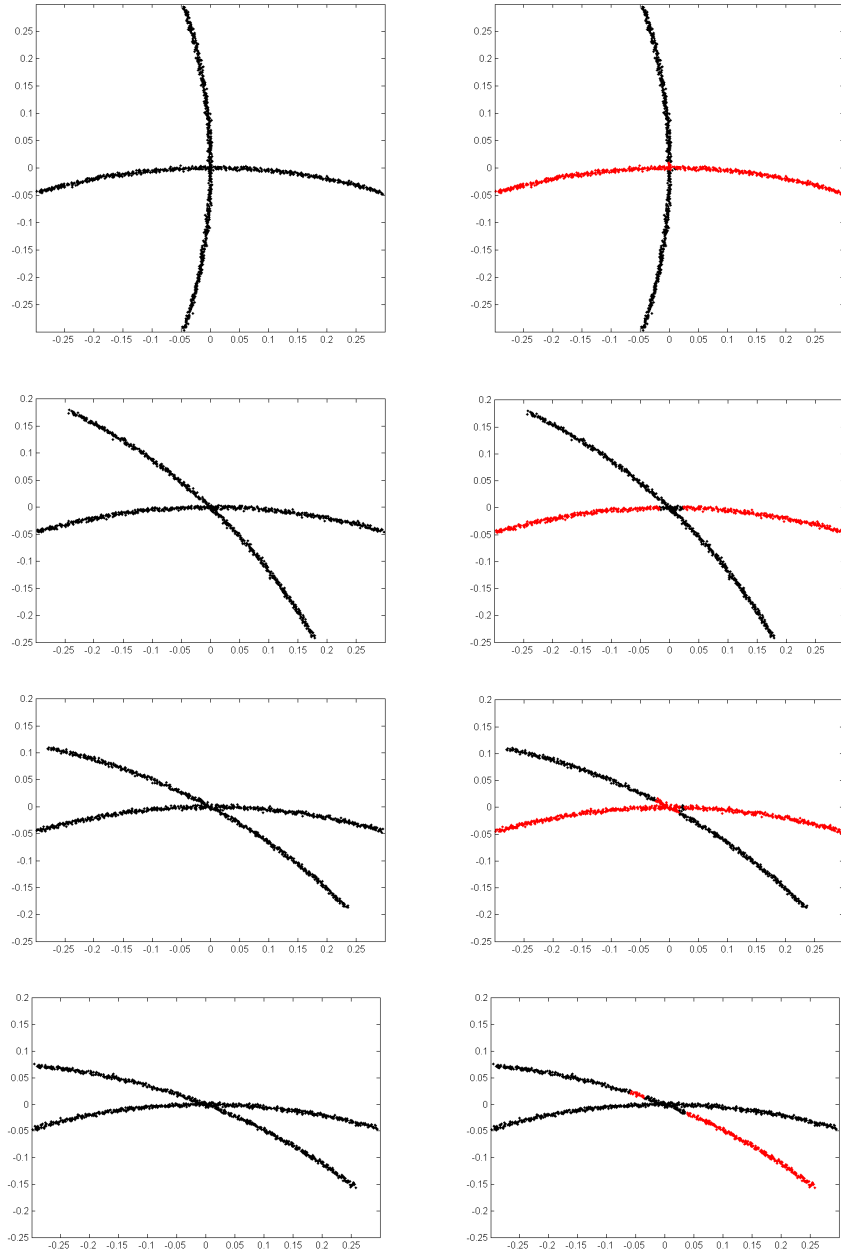


Figure 5: Performance of Algorithm 4 on two curves intersecting at angles  $\frac{\pi}{2}$ ,  $\frac{\pi}{4}$ ,  $\frac{\pi}{6}$ ,  $\frac{\pi}{8}$ .

## 4.2 Comparison with Other Multi-manifold Clustering Algorithms

In this section, we compare our algorithm with several existing algorithms on multi-manifold clustering. While many algorithms have been proposed, we focus on the methods based on spectral clustering, including Sparse Manifold Clustering and Embedding (SMCE) (Elhamifar and Vidal, 2011) and Local Linear Embedding of (Polito and Perona, 2001; Goh and Vidal, 2007). Compared to these methods, a major difference of Algorithm 4 is the size of the affinity matrix  $\mathbf{W}$ : SMCE and LLE each creates an  $n \times n$  affinity matrix while our method creates a smaller affinity matrix of size  $n_0 \times n_0$ , based on the centers chosen in step 0 of Algorithm 4. This difference enables our algorithm to handle large data sets such as  $n > 10^4$ , while these other two methods are computationally expensive due to eigenvalue decomposition of the  $n \times n$  affinity matrix. In order to make a fair comparison, we will run simulations and experiments on small data sets. We modified our algorithm to make it more competitive in such a setting: we modify Steps 0 and 1 in Algorithm 4 slightly and use all data points  $\{\mathbf{x}_i\}_{i=1}^n$  as centers, that is,  $\mathbf{y}_i = \mathbf{x}_i$  for all  $1 \leq i \leq n$ . The motivation is that, when there is no computational constraint on the eigenvalue decomposition (due to small  $n$ ), we may improve Algorithm 4 by constructing a larger affinity matrix by including all points as centers. The resulting algorithm is summarized in Algorithm 5.

---

### Algorithm 5 Spectral Clustering Based on Local PCA (for small data sets)

---

**Input:**

Data points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ ; neighborhood size  $N > 0$ ; spatial scale  $\varepsilon > 0$ , projection scale  $\eta > 0$ ; intrinsic dimension  $d$ ; number of clusters  $K$ .

**Steps:**

**1:** For each  $i = 1, \dots, n$ , compute the sample covariance matrix  $\mathbf{C}_i$  of from the  $N$  nearest neighbors of  $\mathbf{x}_i$ . Let  $\mathbf{Q}_i$  denote the orthogonal projection onto the space spanned by the top  $d$  eigenvectors of  $\mathbf{C}_i$ .

**2:** Compute the following affinities between  $n$  data points:

$$W_{ij} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\varepsilon^2}\right) \cdot \exp\left(-\frac{\|\mathbf{Q}_i - \mathbf{Q}_j\|^2}{\eta^2}\right).$$

**3:** Apply spectral graph partitioning (Algorithm 1) to  $\mathbf{W}$ .

---

First we test the algorithms on a simulated data set of two curves, which is subsampled from the first data set in Figure 5 with 300 data points. We plot the clustering result from Algorithm 5, SMCE, LLE in Figure 6. For Algorithm 5,  $K$  is set to be 10. For SMCE<sup>3</sup>,  $\lambda = 10$  and  $L = 60$ , and we remark that similar results are obtained for a wide range of parameters. For LLE, we follow the implementation in (Polito and Perona, 2001), use 10-nearest neighbors to embed the data set into  $\mathbb{R}^2$  and run  $K$ -means on the embedded data set. It is clear from the figure that Algorithm 5 resolves the problem of intersection well

---

3. In (Elhamifar and Vidal, 2011),  $\lambda$  is the  $\ell_1$ -penalty parameter and at each point the optimization is restricted to the  $L$  nearest neighbors.

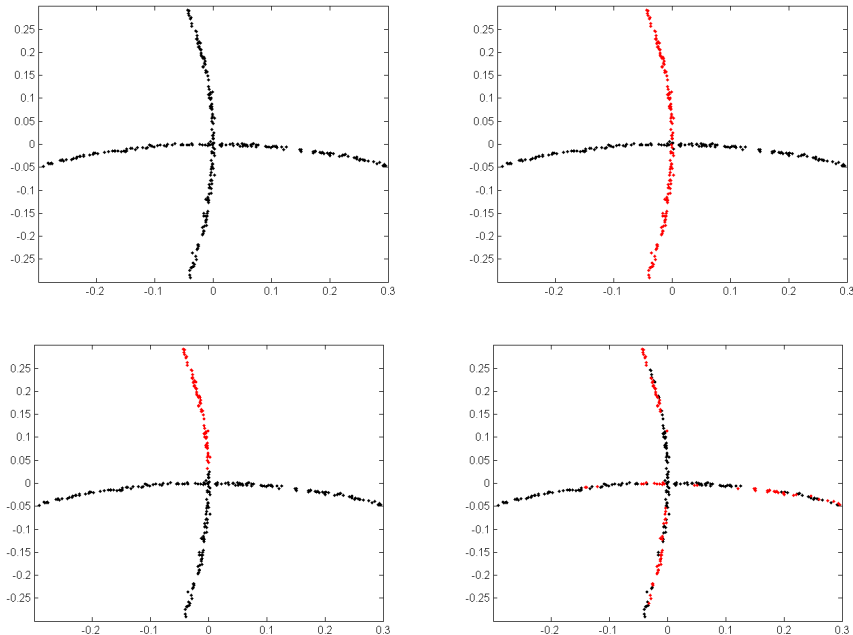


Figure 6: Performance of Algorithm 5 on two intersecting curves. Left top: the input data set. Right top: the clustering result by Algorithm 5. Left bottom: the clustering result by SMCE. Right bottom: the clustering result by LLE.

by using the affinity from estimated local subspaces, while SMCE and LLE tend to give a larger affinity between nearby data points and have difficulties in handling intersection.

Next, we run experiments on the Extended Yale Face Database B (Lee et al., 2005), with the goal of clustering face images of two different subjects. This data set contains face images from 39 subjects, and each subject has 64 images of 192 pixels under varying lightening conditions. In our experiments, we found that the images of a person in this database lie roughly in a 4-dimensional subspace. We preprocess the data set by applying PCA and reducing the dimension to 8. We also “normalize” the covariance of the data set when performing dimension reduction, such that the projected data set has a unit covariance. We record the misclustering rates of Algorithm 5, SMCE and LLE in Table 3. For SMCE, we follow (Elhamifar and Vidal, 2011) by setting  $\lambda = 10$  and we let  $L = 30$ . For Algorithm 5, we let the neighborhood size be 40. From the table, we can see that the two methods perform similarly.

subjects	[1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]	[1,9]
Local PCA	8.59%	11.72%	10.94%	4.69%	8.59%	7.81%	5.47%	7.03%
SMCE	8.59%	11.72%	8.59%	0.00%	4.69%	8.59%	9.38%	4.69%

Table 3: Some misclustering rates for the Extended Yale Face Database B.

## 5. Discussion

We distilled the ideas of Goldberg et al. (2009) and of Kushnir et al. (2006) to cluster points sampled near smooth surfaces. The key ingredient is the use of local PCA to learn about the local spread and orientation of the data, so as to use that information in an affinity when building a neighborhood graph.

In a typical stylized setting for multi-manifold clustering, we established performance bounds for the simple variants described in Algorithm 2 and Algorithm 3, which essentially consist of connecting points that are close in space and orientation, and then extracting the connected components of the resulting graph. Both are shown to resolve general intersections as long as the incidence angle is strictly positive and the parameters are carefully chosen. As is commonly the case in such analyses, our setting can be generalized to other sampling schemes, to multiple intersections, to some features of the surfaces changing with the sample size, and so on, in the spirit of (Arias-Castro et al., 2011; Arias-Castro, 2011; Chen and Lerman, 2009b). We chose to simplify the setup as much as possible while retaining the essential features that make resolving intersecting clusters challenging. The resulting arguments are nevertheless rich enough to satisfy the mathematically thirsty reader. Whether the conditions required in Theorem 1 are optimal in some sense is an interesting and challenging open question for future research. Note that very few optimality results exist for manifold clustering; see (Arias-Castro, 2011) for an example.

We implemented a spectral version of Algorithm 3, described in Algorithm 4, that assumes the intrinsic dimensionality and the number of clusters are known. The resulting approach is very similar to what is offered by Wang et al. (2011) and Gong et al. (2012), although it was developed independently of these works. Algorithm 4 is shown to perform well in some simulated experiments, although it is somewhat sensitive to the choice of parameters. This is the case of all other methods for multi-manifold clustering we know of and choosing the parameters automatically remains an open challenge in the field.

## 6. Proofs

We start with some additional notation. The ambient space is  $\mathbb{R}^D$  unless noted otherwise. For a vector  $\mathbf{v} \in \mathbb{R}^D$ ,  $\|\mathbf{v}\|$  denotes its Euclidean norm and for a real matrix  $\mathbf{M} \in \mathbb{R}^{D \times D}$ ,  $\|\mathbf{M}\|$  denotes the corresponding operator norm. For a point  $\mathbf{x} \in \mathbb{R}^D$  and  $r > 0$ ,  $B(\mathbf{x}, r)$  denotes the open ball of center  $\mathbf{x}$  and radius  $r$ , i.e.,  $B(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^D : \|\mathbf{y} - \mathbf{x}\| < r\}$ . For a set  $S$  and a point  $\mathbf{x}$ , define  $\text{dist}(\mathbf{x}, S) = \inf\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{y} \in S\}$ . For two points  $\mathbf{a}, \mathbf{b}$  in the same Euclidean space,  $\mathbf{b} - \mathbf{a}$  denotes the vector moving  $\mathbf{a}$  to  $\mathbf{b}$ . For a point  $\mathbf{a}$  and a vector  $\mathbf{v}$  in the same Euclidean space,  $\mathbf{a} + \mathbf{v}$  denotes the translate of  $\mathbf{a}$  by  $\mathbf{v}$ . We identify an affine subspace  $T$  with its corresponding linear subspace, for example, when saying that a vector belongs to  $T$ .

For two subspaces  $T$  and  $T'$ , of possibly different dimensions, let  $0 \leq \theta_{\max}(T, T') \leq \pi/2$  denote the largest and by  $\theta_{\min}(T, T')$  the smallest nonzero principal angle between  $T$  and  $T'$  (Stewart and Sun, 1990). When  $\mathbf{v}$  is a vector and  $T$  is a subspace,  $\angle(\mathbf{v}, T) := \theta_{\max}(\mathbb{R}\mathbf{v}, T)$  this is the usual definition of the angle between  $\mathbf{v}$  and  $T$ .

For a subset  $A \subset \mathbb{R}^D$  and positive integer  $d$ ,  $\text{vol}_d(A)$  denotes the  $d$ -dimensional Hausdorff measure of  $A$ , and  $\text{vol}(A)$  is defined as  $\text{vol}_{\dim(A)}(A)$ , where  $\dim(A)$  is the Hausdorff dimension of  $A$ . For a Borel set  $A$ , let  $\lambda_A$  denote the uniform distribution on  $A$ .

For a set  $S \subset \mathbb{R}^D$  with reach at least  $1/\kappa$ , and  $\mathbf{x}$  with  $\text{dist}(\mathbf{x}, S) < 1/\kappa$ , let  $P_S(\mathbf{x})$  denote the metric projection of  $\mathbf{x}$  onto  $S$ , that is, the point on  $S$  closest to  $\mathbf{x}$ . Note that, if  $T$  is an affine subspace, then  $P_T$  is the usual orthogonal projection onto  $T$ , and we let  $\mathbf{P}_T$  denote the orthogonal projection onto the linear subspace of same dimension and parallel to  $T$ . Let  $\mathcal{S}_d(\kappa)$  denote the class of connected,  $C^2$  and compact  $d$ -dimensional submanifolds without boundary embedded in  $\mathbb{R}^D$ , with reach at least  $1/\kappa$ . For a submanifold  $S \in \mathcal{S}_d(\kappa)$ , let  $T_S(\mathbf{x})$  denote the tangent space of  $S$  at  $\mathbf{x} \in S$ .

We will often identify a linear map with its matrix in the canonical basis. For a symmetric (real) matrix  $\mathbf{M}$ , let  $\beta_1(\mathbf{M}) \geq \beta_2(\mathbf{M}) \geq \dots$  denote its eigenvalues in decreasing order.

We say that  $f : \Omega \subset \mathbb{R}^D \rightarrow \mathbb{R}^D$  is  $C$ -Lipschitz if  $\|f(\mathbf{x}) - f(\mathbf{y})\| \leq C\|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \Omega$ .

For two reals  $a$  and  $b$ ,  $a \vee b = \max(a, b)$  and  $a \wedge b = \min(a, b)$ . Additional notation will be introduced as needed.

## 6.1 Preliminaries

This section gathers a number of general results from geometry and probability. We took time to package them into standalone lemmas that could be of potential independent interest, particularly to researchers working in machine learning and computational geometry. When needed, we use  $C$  to denote a constant that may change with each appearance.

### 6.1.1 SMOOTH SURFACES AND THEIR TANGENT SUBSPACES

The following comes directly from (Federer, 1959, Th. 4.18(12)). It gives us a simple criterion for identifying the metric projection of a point on a surface with given reach.

**Lemma 1.** *Consider  $S \in \mathcal{S}_d(\kappa)$  and  $\mathbf{x} \in \mathbb{R}^D$  such that  $\text{dist}(\mathbf{x}, S) < 1/\kappa$ . Then  $\mathbf{s} = P_S(\mathbf{x})$  if and only if  $\|\mathbf{x} - \mathbf{s}\| < 1/\kappa$  and  $\mathbf{x} - \mathbf{s} \perp T_S(\mathbf{s})$ .*

The following result is on approximating a smooth surface near a point by the tangent subspace at that point. It is based on (Federer, 1959, Th. 4.18(2)).

**Lemma 2.** *For  $S \in \mathcal{S}_d(\kappa)$ , and any two points  $\mathbf{s}, \mathbf{s}' \in S$ ,*

$$\text{dist}(\mathbf{s}', T_S(\mathbf{s})) \leq \frac{\kappa}{2} \|\mathbf{s}' - \mathbf{s}\|^2, \tag{10}$$

and when  $\text{dist}(\mathbf{s}', T_S(\mathbf{s})) \leq 1/\kappa$ ,

$$\text{dist}(\mathbf{s}', T_S(\mathbf{s})) \leq \kappa \|P_{T_S(\mathbf{s})}(\mathbf{s}') - \mathbf{s}\|^2. \tag{11}$$

Moreover, for  $\mathbf{t} \in T_S(\mathbf{s})$  such that  $\|\mathbf{s} - \mathbf{t}\| \leq \frac{1}{3\kappa}$ ,

$$\text{dist}(\mathbf{t}, S) \leq \kappa \|\mathbf{t} - \mathbf{s}\|^2. \tag{12}$$

*Proof.* Let  $T$  be short for  $T_S(\mathbf{s})$ . (Federer, 1959, Th. 4.18(2)) says that

$$\text{dist}(\mathbf{s}' - \mathbf{s}, T) \leq \frac{\kappa}{2} \|\mathbf{s}' - \mathbf{s}\|^2. \quad (13)$$

Immediately, we have

$$\text{dist}(\mathbf{s}' - \mathbf{s}, T) = \|\mathbf{s}' - P_T(\mathbf{s}')\| = \text{dist}(\mathbf{s}', T),$$

and (10) comes from that. Based on that and Pythagoras theorem, we have

$$\text{dist}(\mathbf{s}', T) = \|P_T(\mathbf{s}') - \mathbf{s}'\| \leq \frac{\kappa}{2} \|\mathbf{s}' - \mathbf{s}\|^2 = \frac{\kappa}{2} (\|P_T(\mathbf{s}') - \mathbf{s}'\|^2 + \|P_T(\mathbf{s}') - \mathbf{s}\|^2),$$

so that

$$\text{dist}(\mathbf{s}', T) \left(1 - \frac{\kappa}{2} \text{dist}(\mathbf{s}', T)\right) \leq \frac{\kappa}{2} \|P_T(\mathbf{s}') - \mathbf{s}\|^2,$$

and (11) follows easily from that. For (12), let  $r = 1/(3\kappa)$  and  $\mathbf{s}' = P_T^{-1}(\mathbf{t})$ , the latter being well-defined by Lemma 5 below and belongs to  $B(\mathbf{s}, r(1 + \kappa r)) \subset B(\mathbf{s}, 4/(9\kappa))$ . By (10),  $\|\mathbf{s}' - \mathbf{s}\| \leq 8/(81\kappa) < 1/\kappa$ , and by (11),

$$\text{dist}(\mathbf{t}, S) \leq \|\mathbf{t} - \mathbf{s}'\| = \text{dist}(\mathbf{s}', T) \leq \kappa \|\mathbf{t} - \mathbf{s}\|^2.$$

This concludes the proof of (12).  $\square$

We will need a bound on the angle between tangent subspaces on a smooth surface as a function of the distance between the corresponding points of contact.

**Lemma 3** (Boissonnat et al. (2013)). *For  $S \in \mathcal{S}_d(\kappa)$ , and any  $\mathbf{s}, \mathbf{s}' \in S$ ,*

$$\sin \theta_{\max}(T_S(\mathbf{s}), T_S(\mathbf{s}')) \leq 6\kappa \|\mathbf{s}' - \mathbf{s}\|. \quad (14)$$

The following bounds the difference between the metric projection onto a surface and the orthogonal projection onto one of its tangents.

**Lemma 4.** *Consider  $S \in \mathcal{S}_d(\kappa)$  and  $\mathbf{s} \in S$ . Then for any  $\mathbf{x} \in B(S, r)$ , we have*

$$\|P_S(\mathbf{x}) - P_{T_S(\mathbf{s})}(\mathbf{x})\| \leq C_4 \kappa \|\mathbf{x} - \mathbf{s}\|^2,$$

for a numeric constant  $C_4 > 0$ .

*Proof.* Let  $T = T_S(\mathbf{s})$  and define  $\mathbf{t} = P_T(\mathbf{x})$ ,  $\mathbf{s}' = P_S(\mathbf{x})$  and  $\tilde{\mathbf{t}} = P_T(\mathbf{s}')$ , and also  $T' = T_S(\mathbf{s}')$  and  $\theta = \theta_{\max}(T, T')$ . We have

$$\begin{aligned} P_S(\mathbf{x}) &= P_{T'}(\mathbf{x}) = \mathbf{P}_{T'}(\mathbf{x} - \mathbf{s}') + \mathbf{s}' \\ P_T(\mathbf{x}) &= \mathbf{P}_T(\mathbf{x} - \tilde{\mathbf{t}}) + \tilde{\mathbf{t}} = \mathbf{P}_T(\mathbf{x} - \mathbf{s}') + \mathbf{P}_T(\mathbf{s}' - \tilde{\mathbf{t}}) + \tilde{\mathbf{t}}. \end{aligned}$$

Hence,

$$\|P_S(\mathbf{x}) - P_T(\mathbf{x})\| \leq \|\mathbf{P}_{T'}(\mathbf{x} - \mathbf{s}') - \mathbf{P}_T(\mathbf{x} - \mathbf{s}')\| + \|\mathbf{P}_T(\mathbf{s}' - \tilde{\mathbf{t}})\| + \|\mathbf{s}' - \tilde{\mathbf{t}}\|.$$

On the one hand,

$$\|P_T(\mathbf{s}' - \tilde{\mathbf{t}})\| \leq \|\mathbf{s}' - \tilde{\mathbf{t}}\| = \text{dist}(\mathbf{s}', T) \leq \frac{\kappa}{2} \|\mathbf{s}' - \mathbf{s}\|^2 \leq 2\kappa \|\mathbf{x} - \mathbf{s}\|^2,$$

by (10) and the fact that

$$\|\mathbf{s}' - \mathbf{s}\| \leq \|\mathbf{s}' - \mathbf{x}\| + \|\mathbf{x} - \mathbf{s}\| = \text{dist}(\mathbf{x}, S) + \|\mathbf{x} - \mathbf{s}\| \leq 2\|\mathbf{x} - \mathbf{s}\|.$$

On the other hand, applying Lemma 18 (see further down)

$$\|P_{T'}(\mathbf{x} - \mathbf{s}') - P_T(\mathbf{x} - \mathbf{s}')\| \leq \|P_{T'} - P_T\| \|\mathbf{x} - \mathbf{s}'\| = (\sin \theta) \|\mathbf{x} - \mathbf{s}'\|,$$

with  $\|\mathbf{x} - \mathbf{s}'\| \leq \|\mathbf{x} - \mathbf{s}\|$  and, applying Lemma 3,

$$\sin \theta \leq 6\kappa \|\mathbf{s}' - \mathbf{s}\| \leq 12\kappa \|\mathbf{x} - \mathbf{s}\|.$$

All together, we conclude.  $\square$

Below we state some properties of a projection onto a tangent subspace. A result similar to the first part was proved in (Arias-Castro et al., 2011, Lem. 2) based on results in (Niyogi et al., 2008), but the arguments are simpler here and the constants are sharper.

**Lemma 5.** *There is a numeric constant  $C_5 \geq 1$  such that the following holds. Take  $S \in \mathcal{S}_d(\kappa)$ ,  $\mathbf{s} \in S$  and  $r \leq 1/C_5\kappa$ , and let  $T$  be short for  $T_S(\mathbf{s})$ .  $P_T$  is injective on  $B(\mathbf{s}, r) \cap S$  and its image contains  $B(\mathbf{s}, r') \cap T$ , where  $r' := (1 - C_5(\kappa r)^2)r$ . Moreover,  $P_T^{-1}$  has Lipschitz constant bounded by  $1 + C_5(\kappa r)^2 \leq 1 + \kappa r$  over  $B(\mathbf{s}, r) \cap T$ .*

*Proof.* Take  $\mathbf{s}', \mathbf{s}'' \in S$  distinct such that  $P_T(\mathbf{s}') = P_T(\mathbf{s}'')$ . Equivalently,  $\mathbf{s}'' - \mathbf{s}'$  is perpendicular to  $T$ . Let  $T'$  be short for  $T_S(\mathbf{s}')$ . By (13) and the fact that  $\text{dist}(\mathbf{v}, T) = \|\mathbf{v}\| \sin \angle(\mathbf{v}, T)$  for any vector  $\mathbf{v}$  and any linear subspace  $T$ , we have

$$\sin \angle(\mathbf{s}'' - \mathbf{s}', T') \leq \frac{\kappa}{2} \|\mathbf{s}'' - \mathbf{s}'\|,$$

and by (14),

$$\sin \theta_{\max}(T, T') \leq 6\kappa \|\mathbf{s} - \mathbf{s}'\|.$$

Now, by the triangle inequality,

$$\frac{\pi}{2} = \angle(\mathbf{s}'' - \mathbf{s}', T) \leq \angle(\mathbf{s}'' - \mathbf{s}', T') + \theta_{\max}(T, T'),$$

so that

$$\frac{\kappa}{2} \|\mathbf{s}'' - \mathbf{s}'\| \wedge 1 \geq \frac{\pi}{2} - \sin^{-1}(6\kappa \|\mathbf{s}' - \mathbf{s}\| \wedge 1).$$

When  $\|\mathbf{s}' - \mathbf{s}\| \leq 1/12\kappa$ , the RHS is bounded from below by  $\pi/2 - \sin^{-1}(1/2)$ , which then implies that  $\frac{\kappa}{2} \|\mathbf{s}'' - \mathbf{s}'\| \geq \sin(\pi/2 - \sin^{-1}(1/2)) = \sqrt{3}/2$ , that is,  $\|\mathbf{s}'' - \mathbf{s}'\| \geq \sqrt{3}/\kappa$ . This precludes the situation where  $\mathbf{s}', \mathbf{s}'' \in B(\mathbf{s}, 1/12\kappa)$ , so that  $P_T$  is injective on  $B(\mathbf{s}, r)$  when  $r \leq 1/12\kappa$ .

The same arguments imply that  $P_T$  is an open map on  $R := B(\mathbf{s}, r) \cap S$ . In particular,  $P_T(R)$  contains an open ball in  $T$  centered at  $\mathbf{s}$  and  $P_T(\partial R) = \partial P_T(R)$ , with  $\partial R = S \cap$

$\partial B(\mathbf{s}, r)$  since  $\partial S = \emptyset$ . Now take any ray out of  $\mathbf{s}$  within  $T$ , which is necessarily of the form  $\mathbf{s} + \mathbb{R}_+ \mathbf{v}$ , where  $\mathbf{v}$  is a unit vector in  $T$ . Let  $\mathbf{t}_a = \mathbf{s} + a\mathbf{v} \in T$  for  $a \in [0, \infty)$ . Let  $a_*$  be the infimum over all  $a > 0$  such that  $\mathbf{t}_a \in P_T(R)$ . Note that  $a_* > 0$  and  $\mathbf{t}_{a_*} \in P_T(\partial R)$ , so that there is  $\mathbf{s}_* \in \partial R$  such that  $P_T(\mathbf{s}_*) = \mathbf{t}_{a_*}$ . Let  $\mathbf{s}_a = P_T^{-1}(\mathbf{t}_a)$ , which is well-defined on  $[0, a_*]$  by definition of  $a_*$  and the fact that  $P_T$  is injective on  $R$ . Let  $J_{\mathbf{t}}$  denote the differential of  $P_T^{-1}$  at  $\mathbf{t}$ . We have that  $\dot{\mathbf{s}}_a = J_{\mathbf{t}_a} \mathbf{v}$  is the unique vector in  $T_a := T_S(\mathbf{s}_a)$  such that  $P_T(\dot{\mathbf{s}}_a) = \mathbf{v}$ . Elementary geometry shows that

$$\|P_T(\dot{\mathbf{s}}_a)\| = \|\dot{\mathbf{s}}_a\| \cos \angle(\dot{\mathbf{s}}_a, T) \geq \|\dot{\mathbf{s}}_a\| \cos \theta_{\max}(T_a, T),$$

with

$$\cos \theta_{\max}(T_a, T) \geq \cos [\sin^{-1}(6\kappa\|\mathbf{s}_a - \mathbf{s}\| \wedge 1)] \geq \zeta := 1 - (6\kappa r)^2,$$

by (14) and fact that  $\|\mathbf{s}_a - \mathbf{s}\| \leq r$  (and assuming  $6\kappa r \leq 1$ ). Since  $\|P_T(\dot{\mathbf{s}}_a)\| = \|\mathbf{v}\| = 1$ , we have  $\|\dot{\mathbf{s}}_a\| \leq 1/\zeta$ , and this holds for all  $a < a_*$ . So we can extend  $\mathbf{s}_a$  to  $[0, a_*]$  into a Lipschitz function with constant  $1/\zeta$ . Together with the fact that  $\mathbf{s}_* \in \partial B(\mathbf{s}, r)$ , this implies that

$$r = \|\mathbf{s}_* - \mathbf{s}\| = \|\mathbf{s}_{a_*} - \mathbf{s}_0\| \leq a_*/\zeta.$$

Hence,  $a_* \geq \zeta r$  and therefore  $P_T(R)$  contains  $B(\mathbf{s}, \zeta r) \cap T$  as stated.

For the last part, assume  $r \leq 1/C\kappa$ , with  $C$  large enough that there is a unique  $h \leq 1/12\kappa$  such that  $\zeta h = r$ , where  $\zeta$  is redefined as  $\zeta := 1 - (6\kappa h)^2$ . Take  $\mathbf{t}' \in B(\mathbf{s}, r) \cap T$  and let  $\mathbf{s}' = P_T^{-1}(\mathbf{t}')$  and  $T' = T_S(\mathbf{s}')$ . We saw that  $P_T^{-1}$  is Lipschitz with constant  $1/\zeta$  on any ray emanating from  $\mathbf{s}$  of length  $\zeta h = r$ , so that  $\|\mathbf{s}' - \mathbf{s}\| \leq (1/\zeta)\|\mathbf{t}' - \mathbf{s}\| \leq r/\zeta = h$ . The differential of  $P_T$  at  $\mathbf{s}'$  is  $P_T$  itself, seen as a linear map between  $T'$  and  $T$ . Then for any vector  $\mathbf{u} \in T'$ , we have

$$\|P_T(\mathbf{u})\| = \|\mathbf{u}\| \cos \angle(\mathbf{u}, T) \geq \|\mathbf{u}\| \cos \theta_{\max}(T', T),$$

with

$$\cos \theta_{\max}(T', T) \geq \cos [\sin^{-1}(6\kappa\|\mathbf{s}' - \mathbf{s}\|)] \geq 1 - (6\kappa h)^2 = \zeta,$$

as before. Hence,  $\|J_{\mathbf{t}'}\| \leq 1/\zeta$ , and we proved this for all  $\mathbf{t}' \in B(\mathbf{s}, r) \cap T$ . This last set being convex, we can apply Taylor's theorem and get that  $P_T^{-1}$  is Lipschitz on that set with constant  $1/\zeta$ . We then note that  $\zeta = 1 + O(\kappa r)^2$ .  $\square$

### 6.1.2 VOLUMES AND UNIFORM DISTRIBUTIONS

Below is a result that quantifies how much the volume of a set changes when applying a Lipschitz map. This is well-known in measure theory and we only provide a proof for completeness.

**Lemma 6.** *Suppose  $\Omega$  is a measurable subset of  $\mathbb{R}^D$  and  $f : \Omega \subset \mathbb{R}^D \rightarrow \mathbb{R}^D$  is  $C$ -Lipschitz. Then for any measurable set  $A \subset \Omega$  and real  $d > 0$ ,  $\text{vol}_d(f(A)) \leq C^d \text{vol}_d(A)$ .*

*Proof.* By definition,

$$\text{vol}_d(A) = \lim_{t \rightarrow 0} V_d^t(A), \quad V_d^t(A) := \inf_{(R_i) \in \mathcal{R}^t(A)} \sum_{i \in \mathbb{N}} \text{diam}(R_i)^d,$$



where  $\mathcal{R}^t(A)$  is the class of countable sequences  $(R_i : i \in \mathbb{N})$  of subsets of  $\mathbb{R}^D$  such that  $A \subset \bigcup_i R_i$  and  $\text{diam}(R_i) < t$  for all  $i$ . Since  $f$  is  $C$ -Lipschitz,  $\text{diam}(f(R)) \leq C \text{diam}(R)$  for any  $R \subset \Omega$ . Hence, for any  $(R_i) \in \mathcal{R}^t(A)$ ,  $(f(R_i)) \in \mathcal{R}^{Ct}(f(A))$ . This implies that

$$V_d^{Ct}(f(A)) \leq \sum_{i \in \mathbb{N}} \text{diam}(f(R_i))^d \leq C^d \sum_{i \in \mathbb{N}} \text{diam}(R_i)^d.$$

Taking the infimum over  $(R_i) \in \mathcal{R}^t(A)$ , we get  $V_d^{Ct}(f(A)) \leq C^d V_d^t(A)$ , and we conclude by taking the limit as  $t \rightarrow 0$ , noticing that  $V_d^{Ct}(f(A)) \rightarrow \text{vol}_d(f(A))$  while  $V_d^t(A) \rightarrow \text{vol}_d(A)$ .  $\square$

We compare below two uniform distributions. For two Borel probability measures  $P$  and  $Q$  on  $\mathbb{R}^D$ ,  $\text{TV}(P, Q)$  denotes their total variation distance, meaning,

$$\text{TV}(P, Q) = \sup\{|P(A) - Q(A)| : A \text{ Borel}\}.$$

Remember that for a Borel set  $A$ ,  $\lambda_A$  denotes the uniform distribution on  $A$ .

**Lemma 7.** *Suppose  $A$  and  $B$  are two Borel subsets of  $\mathbb{R}^D$ . Then*

$$\text{TV}(\lambda_A, \lambda_B) \leq 4 \frac{\text{vol}(A \triangle B)}{\text{vol}(A \cup B)}.$$

*Proof.* If  $A$  and  $B$  are not of same dimension, say  $\dim(A) > \dim(B)$ , then  $\text{TV}(\lambda_A, \lambda_B) = 1$  since  $\lambda_A(B) = 0$  while  $\lambda_B(B) = 1$ . And we also have

$$\text{vol}(A \triangle B) = \text{vol}_{\dim(A)}(A \triangle B) = \text{vol}_{\dim(A)}(A) = \text{vol}(A),$$

and

$$\text{vol}(A \cup B) = \text{vol}_{\dim(A)}(A \cup B) = \text{vol}_{\dim(A)}(A) = \text{vol}(A),$$

in both cases because  $\text{vol}_{\dim(A)}(B) = 0$ . So the result works in that case.

Therefore assume that  $A$  and  $B$  are of same dimension. Assume WLOG that  $\text{vol}(A) \geq \text{vol}(B)$ . For any Borel set  $U$ ,

$$\lambda_A(U) - \lambda_B(U) = \frac{\text{vol}(A \cap U)}{\text{vol}(A)} - \frac{\text{vol}(B \cap U)}{\text{vol}(B)},$$

so that

$$\begin{aligned} |\lambda_A(U) - \lambda_B(U)| &= \left| \frac{\text{vol}(A \cap U)}{\text{vol}(A)} - \frac{\text{vol}(B \cap U)}{\text{vol}(A)} + \frac{\text{vol}(B \cap U)}{\text{vol}(A)} - \frac{\text{vol}(B \cap U)}{\text{vol}(B)} \right| \\ &\leq \frac{|\text{vol}(A \cap U) - \text{vol}(B \cap U)|}{\text{vol}(A)} + \frac{\text{vol}(B \cap U)}{\text{vol}(B)} \frac{|\text{vol}(A) - \text{vol}(B)|}{\text{vol}(A)} \\ &\leq \frac{2 \text{vol}(A \triangle B)}{\text{vol}(A)}, \end{aligned}$$

and we conclude with the fact that  $\text{vol}(A \cup B) \leq \text{vol}(A) + \text{vol}(B) \leq 2 \text{vol}(A)$ .  $\square$

We now look at the projection of the uniform distribution on a neighborhood of a surface onto a tangent subspace. For a Borel probability measure  $P$  and measurable function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^D$ ,  $P^f$  denotes the push-forward (Borel) measure defined by  $P^f(A) = P(f^{-1}(A))$ .

**Lemma 8.** *Suppose  $A \subset \mathbb{R}^D$  is Borel and  $f : A \rightarrow \mathbb{R}^D$  is invertible on  $f(A)$ , and that both  $f$  and  $f^{-1}$  are  $C$ -Lipschitz. Then*

$$\text{TV}(\lambda_A^f, \lambda_{f(A)}) \leq 2(C^{\dim(A)} - 1).$$

*Proof.* First, note that  $A$  and  $f(A)$  are both of same dimension, and that  $C \geq 1$  necessarily. Let  $d$  be short for  $\dim(A)$ . Take  $U \subset f(A)$  Borel and let  $V = f^{-1}(U)$ . Then

$$\lambda_A^f(U) = \frac{\text{vol}(A \cap V)}{\text{vol}(A)}, \quad \lambda_{f(A)}(U) = \frac{\text{vol}(f(A) \cap U)}{\text{vol}(f(A))},$$

and as in (16),

$$|\lambda_A^f(U) - \lambda_{f(A)}(U)| \leq \frac{|\text{vol}(A \cap V) - \text{vol}(f(A) \cap U)|}{\text{vol}(A)} + \frac{|\text{vol}(A) - \text{vol}(f(A))|}{\text{vol}(A)}.$$

$f$  being invertible, we have  $f(A \cap V) = f(A) \cap U$  and  $f^{-1}(f(A) \cap U) = A \cap V$ . Therefore, applying Lemma 6, we get

$$C^{-d} \leq \frac{\text{vol}(f(A) \cap U)}{\text{vol}(A \cap V)} \leq C^d,$$

so that

$$|\text{vol}(A \cap V) - \text{vol}(f(A) \cap U)| \leq (C^d - 1) \text{vol}(A \cap V) \leq (C^d - 1) \text{vol}(A).$$

And taking  $V = A$ , we also get

$$|\text{vol}(A) - \text{vol}(f(A))| \leq (C^d - 1) \text{vol}(A).$$

From this we conclude. □

Now comes a technical result on the intersection of a smooth surface and a ball.

**Lemma 9.** *There is a constant  $C_9 \geq 3$  depending only on  $d$  such that the following is true. Take  $S \in \mathcal{S}_d(\kappa)$ ,  $r < \frac{1}{C_9\kappa}$  and  $\mathbf{x} \in \mathbb{R}^D$  such that  $\text{dist}(\mathbf{x}, S) < \kappa$ . Let  $\mathbf{s} = P_S(\mathbf{x})$  and  $T = T_S(\mathbf{s})$ . Then*

$$\text{vol}(P_T(S \cap B(\mathbf{x}, r)) \triangle (T \cap B(\mathbf{x}, r))) \leq C_9\kappa(\|\mathbf{x} - \mathbf{s}\| + \kappa r^2) \text{vol}(T \cap B(\mathbf{x}, r)).$$

*Proof.* If  $\text{dist}(\mathbf{x}, S) > r$ , then  $S \cap B(\mathbf{x}, r) = T \cap B(\mathbf{x}, r) = \emptyset$ , and if  $\text{dist}(\mathbf{x}, S) = r$ , then  $S \cap B(\mathbf{x}, r) = T \cap B(\mathbf{x}, r) = \{\mathbf{s}\}$ , and in both cases the inequality holds trivially. So it suffices to consider the case where  $\text{dist}(\mathbf{x}, S) < r$ .

Let  $A_r = B(\mathbf{s}, r)$ ,  $B_r = B(\mathbf{x}, r)$  and  $g = P_T$  for short. Note that  $T \cap B_r = T \cap A_{r_0}$  where  $r_0 := (r^2 - \delta^2)^{1/2}$  and  $\delta := \|\mathbf{x} - \mathbf{s}\|$ . Take  $\mathbf{s}_1 \in S \cap B_r$  such that  $g(\mathbf{s}_1)$  is farthest from  $\mathbf{s}$ , so that  $g(S \cap B_r) \subset A_{r_1}$  where  $r_1 := \|\mathbf{s} - g(\mathbf{s}_1)\|$ —note that  $r_1 \leq r$ . Let  $\ell_1 = \|\mathbf{s}_1 - g(\mathbf{s}_1)\|$  and  $\mathbf{y}_1$  be the orthogonal projection of  $\mathbf{s}_1$  onto the line  $(\mathbf{x}, \mathbf{s})$ . By Pythagoras theorem,

we have  $\|\mathbf{x} - \mathbf{s}_1\|^2 = \|\mathbf{x} - \mathbf{y}_1\|^2 + \|\mathbf{y}_1 - \mathbf{s}_1\|^2$ . We have  $\|\mathbf{x} - \mathbf{s}_1\| \leq r$  and  $\|\mathbf{y}_1 - \mathbf{s}_1\| = \|\mathbf{s} - g(\mathbf{s}_1)\| = r_1$ . And because  $\ell_1 \leq \kappa r_1^2 < r$  by (11), either  $\mathbf{y}_1$  is between  $\mathbf{x}$  and  $\mathbf{s}$ , in which case  $\|\mathbf{x} - \mathbf{y}_1\| = \delta - \ell_1$ , or  $\mathbf{s}$  is between  $\mathbf{x}$  and  $\mathbf{y}_1$ , in which case  $\|\mathbf{x} - \mathbf{y}_1\| = \delta + \ell_1$ . In any case,  $r^2 \geq r_1^2 + (\delta - \ell_1)^2$ , which together with  $\ell_1 \leq \kappa r_1^2$  implies  $r_1^2 \leq r^2 - \delta^2 + 2\delta\ell_1 \leq r_0^2 + 2\kappa r_1^2 \delta$ , leading to  $r_1 \leq (1 - 2\kappa\delta)^{-1/2} r_0 \leq (1 + 4\kappa\delta)r_0$  after noticing that  $\delta \leq r < 1/(3\kappa)$ . From  $g(S \cap B_r) \subset T \cap A_{r_1}$ , we get

$$\begin{aligned} \text{vol}(g(S \cap B_r) \setminus (T \cap B_r)) &\leq \text{vol}(T \cap A_{r_1}) - \text{vol}(T \cap A_{r_0}) \\ &= ((r_1/r_0)^d - 1) \text{vol}(T \cap A_{r_0}). \end{aligned}$$

We follow similar arguments to get a sort of reverse relationship. Take  $\mathbf{s}_2 \in S \cap B_r$  such that  $g(S \cap B_r) \supset T \cap A_{r_2}$ , where  $r_2 := \|\mathbf{s} - g(\mathbf{s}_2)\|$  is largest. Assuming  $r$  is small enough, by Lemma 5,  $g^{-1}$  is well-defined on  $T \cap A_r$ , so that necessarily  $\mathbf{s}_2 \in \partial B_r$ . Let  $\ell_2 = \|\mathbf{s}_2 - g(\mathbf{s}_2)\|$  and  $\mathbf{y}_2$  be the orthogonal projection of  $\mathbf{s}_2$  onto the line  $(\mathbf{x}, \mathbf{s})$ . By Pythagoras theorem, we have  $\|\mathbf{x} - \mathbf{s}_2\|^2 = \|\mathbf{x} - \mathbf{y}_2\|^2 + \|\mathbf{y}_2 - \mathbf{s}_2\|^2$ . We have  $\|\mathbf{x} - \mathbf{s}_2\| = r$  and  $\|\mathbf{y}_2 - \mathbf{s}_2\| = \|\mathbf{s} - g(\mathbf{s}_2)\| = r_2$ . And by the triangle inequality,  $\|\mathbf{x} - \mathbf{y}_2\| \leq \|\mathbf{x} - \mathbf{s}\| + \|\mathbf{y}_2 - \mathbf{s}\| = \delta + \ell_2$ . Hence,  $r^2 \leq r_2^2 + (\delta + \ell_2)^2$ , which together with  $\ell_2 \leq \kappa r_2^2$  by (11), implies  $r_2^2 \geq r^2 - \delta^2 - 2\delta\ell_2 - \ell_2^2 \geq r_0^2 - (2\delta + \kappa r^2)\kappa r_2^2$ , leading to  $r_2 \geq (1 + 2\kappa\delta + \kappa^2 r^2)^{-1/2} r_0 \geq (1 - 2\kappa\delta - \kappa^2 r^2)r_0$ . From  $g(S \cap B_r) \supset T \cap A_{r_2}$ , we get

$$\begin{aligned} \text{vol}((T \cap B_r) \setminus g(S \cap B_r)) &\leq \text{vol}(T \cap A_{r_0}) - \text{vol}(T \cap A_{r_2}) \\ &= (1 - (r_2/r_0)^d) \text{vol}(T \cap A_{r_0}). \end{aligned}$$

All together, we have

$$\begin{aligned} \text{vol}(g(S \cap B_r) \Delta (T \cap B_r)) &\leq ((r_1/r_0)^d - (r_2/r_0)^d) \text{vol}(T \cap A_{r_0}) \\ &\leq ((1 + 4\kappa\delta)^d - (1 - 2\kappa\delta - \kappa^2 r^2)^d) \text{vol}((T \cap B_r)) \\ &\leq C d \kappa (\delta + \kappa r^2) \text{vol}((T \cap B_r)), \end{aligned}$$

when  $\delta \leq r \leq 1/(C\kappa)$  and  $C$  is a large-enough numerical constant.  $\square$

We bound below the  $d$ -volume of a the intersection of a ball with a smooth surface. Although it could be obtained as a special case of Lemma 9, we provide a direct proof because this result is at the cornerstone of many results in the literature on sampling points uniformly on a smooth surface.

**Lemma 10.** *Suppose  $S \in \mathcal{S}_d(\kappa)$ . Then for any  $\mathbf{s} \in S$  and  $r < \frac{1}{(d\sqrt{C_5})\kappa}$ , we have*

$$1 - 2d\kappa r \leq \frac{\text{vol}(S \cap B(\mathbf{s}, r))}{\text{vol}(T \cap B(\mathbf{s}, r))} \leq 1 + 2d\kappa r,$$

where  $T := T_S(\mathbf{s})$  is the tangent subspace of  $S$  at  $\mathbf{s}$ .

*Proof.* Let  $T = T_S(\mathbf{s})$ ,  $B_r = B(\mathbf{s}, r)$  and  $g = P_T$  for short. By Lemma 5,  $g$  is 1-Lipschitz and  $g^{-1}$  is  $(1 + \kappa r)$ -Lipschitz on  $T \cap B_r$ , so by Lemma 6 we have

$$(1 + \kappa r)^{-d} \leq \frac{\text{vol}(g(S \cap B_r))}{\text{vol}(S \cap B_r)} \leq 1.$$

That  $g^{-1}$  is Lipschitz with constant  $1 + \kappa r$  on  $g(S \cap B_r)$  also implies that  $g(S \cap B_r)$  contains  $T \cap B_{r'}$  where  $r' := r/(1 + \kappa r)$ . From this, and the fact that  $g(S \cap B_r) \subset T \cap B_r$ , we get

$$1 \leq \frac{\text{vol}(T \cap B_r)}{\text{vol}(g(S \cap B_r))} \leq \frac{\text{vol}(T \cap B_r)}{\text{vol}(T \cap B_{r'})} = \frac{r^d}{r'^d} = (1 + \kappa r)^d.$$

We therefore have

$$\text{vol}(S \cap B_r) \geq \text{vol}(g(S \cap B_r)) \geq (1 + \kappa r)^{-d} \text{vol}(T \cap B_r),$$

and

$$\text{vol}(S \cap B_r) \leq (1 + \kappa r)^d \text{vol}(g(S \cap B_r)) \leq (1 + \kappa r)^d \text{vol}(T \cap B_r).$$

And we conclude with the inequality  $(1 + x)^d \leq 1 + 2dx$  valid for any  $x \in [0, 1/d]$  and any  $d \geq 1$ .  $\square$

We now look at the density of a sample from the uniform on a smooth, compact surface.

**Lemma 11.** *Consider  $S \in \mathcal{S}_d(\kappa)$ . There is a constant  $C_{11} > 0$  depending on  $S$  such that the following is true. Sample  $n$  points  $\mathbf{s}_1, \dots, \mathbf{s}_n$  independently and uniformly at random from  $S$ . Take  $0 < r < 1/(C_{11}\kappa)$ . Then with probability at least  $1 - C_{11}r^{-d} \exp(-nr^d/C_{11})$ , any ball of radius  $r$  with center on  $S$  has between  $nr^d/C_{11}$  and  $C_{11}nr^d$  sample points.*

*Proof.* For a set  $R$ , let  $N(R)$  denote the number of sample points in  $R$ . For any  $R$  measurable,  $N(R) \sim \text{Bin}(n, p_R)$ , where  $p_R := \text{vol}(R \cap S)/\text{vol}(S)$ . Let  $\mathbf{x}_1, \dots, \mathbf{x}_m$  be a maximal  $(r/2)$ -packing of  $S$ , and let  $B_i = B(\mathbf{x}_i, r/4) \cap S$ . For any  $\mathbf{s} \in S$ , there is  $i$  such that  $\|\mathbf{s} - \mathbf{x}_i\| \leq r/2$ , which implies  $B_i \subset B(\mathbf{s}, r)$  by the triangle inequality. Hence,  $\min_{\mathbf{s} \in S} N(B(\mathbf{s}, r)) \geq \min_i N(B_i)$ .

By the fact that  $B_i \cap B_j = \emptyset$  for  $i \neq j$ ,

$$\text{vol}(S) \geq \sum_{i=1}^m \text{vol}(B_i) \geq m \min_i \text{vol}(B_i),$$

and assuming that  $\kappa r$  is small enough, we have

$$\min_i \text{vol}(B_i) \geq \frac{\omega_d}{2} (r/4)^d,$$

by Lemma 10, where  $\omega_d$  is the volume of the  $d$ -dimensional unit ball. This leads to  $m \leq Cr^{-d}$  and  $p := \min_i p_{B_i} \geq r^d/C$ , when  $\kappa r \leq 1/C$ , where  $C > 0$  depends only on  $S$ .

Now, applying Bernstein's inequality to the binomial distribution, we get

$$\mathbb{P}(N(B_i) \leq np/2) \leq \mathbb{P}(N(B_i) \leq np_{B_i}/2) \leq e^{-(3/32)np_{B_i}} \leq e^{-(3/32)np}.$$

We follow this with the union bound, to get

$$\mathbb{P}\left(\min_{\mathbf{s} \in S} N(B(\mathbf{s}, r)) \leq nr^d/(2C)\right) \leq me^{-(3/32)np} \leq Cr^{-d} e^{-\frac{3}{32C}nr^d}.$$

From this the lower bound follows. The proof of the upper bound is similar.  $\square$

Next, we bound the volume of the symmetric difference between two balls.

**Lemma 12.** *Take  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  and  $0 < \delta \leq 1$ . Then*

$$\frac{\text{vol}(B(\mathbf{x}, \delta) \triangle B(\mathbf{y}, 1))}{2 \text{vol}(B(0, 1))} \leq 1 - (1 - \|\mathbf{x} - \mathbf{y}\|)_+^d \wedge \delta^d.$$

*Proof.* It suffices to prove the result when  $\|\mathbf{x} - \mathbf{y}\| < 1$ . In that case, with  $\gamma := (1 - \|\mathbf{x} - \mathbf{y}\|) \wedge \delta$ , we have  $B(\mathbf{x}, \gamma) \subset B(\mathbf{x}, \delta) \cap B(\mathbf{y}, 1)$ , so that

$$\begin{aligned} \text{vol}(B(\mathbf{x}, \delta) \triangle B(\mathbf{y}, 1)) &= \text{vol}(B(\mathbf{x}, \delta)) + \text{vol}(B(\mathbf{y}, 1)) - 2 \text{vol}(B(\mathbf{x}, \delta) \cap B(\mathbf{y}, 1)) \\ &\leq 2 \text{vol}(B(\mathbf{y}, 1)) - 2 \text{vol}(B(\mathbf{x}, \gamma)) \\ &= 2 \text{vol}(B(0, 1))(1 - \gamma^d). \end{aligned}$$

This concludes the proof. □

### 6.1.3 COVARIANCES

The result below describes explicitly the covariance matrix of the uniform distribution over the unit ball of a subspace.

**Lemma 13.** *Let  $T$  be a subspace of dimension  $d$ . Then the covariance matrix of the uniform distribution on  $T \cap B(0, a)$  (seen as a linear map) is equal to  $a^2 c \mathbf{P}_T$ , where  $c := \frac{1}{d+2}$ .*

*Proof.* Assume WLOG that  $T = \mathbb{R}^d \times \{0\}$ . Let  $X$  be distributed according to the uniform distribution on  $T \cap B(0, 1)$  and let  $R = \|X\|$ . Note that

$$\mathbb{P}(R \leq r) = \frac{\text{vol}(T \cap B(0, r))}{\text{vol}(T \cap B(0, 1))} = r^d, \quad \forall r \in [0, 1].$$

By symmetry,  $\mathbb{E}(X_i X_j) = 0$  if  $i \neq j$ , while

$$\mathbb{E}(X_1^2) = \frac{1}{d} \mathbb{E}(X_1^2 + \dots + X_d^2) = \frac{1}{d} \mathbb{E}(R^2) = \frac{1}{d} \int_0^1 r^2 (dr^{d-1}) dr = \frac{1}{d+2}.$$

Now the covariance matrix of the uniform distribution on  $T \cap B(0, a)$  equals  $\mathbb{E}[(aX)(aX)^\top] = a^2 \mathbb{E}[XX^\top]$ , which is exactly the representation of  $a^2 c \mathbf{P}_T$  in the canonical basis of  $\mathbb{R}^D$ . □

We now show that a bound on the total variation distance between two compactly supported distributions implies a bound on the difference between their covariance matrices. For a measure  $P$  on  $\mathbb{R}^D$  and an integrable function  $f$ , let  $P(f)$  denote the integral of  $f$  with respect to  $P$ , that is,

$$P(f) = \int f(x) P(dx),$$

and let  $\mathbb{E}(P) = P(\mathbf{x})$  and  $\text{Cov}(P) = P(\mathbf{x}\mathbf{x}^\top) - P(\mathbf{x})P(\mathbf{x})^\top$  denote the mean and covariance matrix of  $P$ , respectively.

**Lemma 14.** *Suppose  $\lambda$  and  $\nu$  are two Borel probability measures on  $\mathbb{R}^d$  supported on  $B(0, 1)$ . Then*

$$\|\mathbb{E}(\lambda) - \mathbb{E}(\nu)\| \leq \sqrt{d} \text{TV}(\lambda, \nu), \quad \|\text{Cov}(\lambda) - \text{Cov}(\nu)\| \leq 3d \text{TV}(\lambda, \nu).$$

*Proof.* Let  $f_k(\mathbf{t}) = t_k$  when  $\mathbf{t} = (t_1, \dots, t_d)$ , and note that  $|f_k(\mathbf{t})| \leq 1$  for all  $k$  and all  $\mathbf{t} \in B(0, 1)$ . By the fact that

$$\text{TV}(\lambda, \nu) = \sup\{\lambda(f) - \nu(f) : f : \mathbb{R}^d \rightarrow \mathbb{R} \text{ measurable with } |f| \leq 1\},$$

we have

$$|\lambda(f_k) - \nu(f_k)| \leq \text{TV}(\lambda, \nu), \quad \forall k = 1, \dots, d.$$

Therefore,

$$\|\mathbb{E}(\lambda) - \mathbb{E}(\nu)\|^2 = \sum_{k=1}^d (\lambda(f_k) - \nu(f_k))^2 \leq d \text{TV}(\lambda, \nu)^2,$$

which proves the first part.

Similarly, let  $f_{k\ell}(\mathbf{t}) = t_k t_\ell$ . Since  $|f_{k\ell}(\mathbf{t})| \leq 1$  for all  $k, \ell$  and all  $\mathbf{t} \in B(0, 1)$ , we have

$$|\lambda(f_{k\ell}) - \nu(f_{k\ell})| \leq \text{TV}(\lambda, \nu), \quad \forall k, \ell = 1, \dots, d.$$

Since for any probability measure  $\mu$  on  $\mathbb{R}^d$ ,

$$\text{Cov}(\mu) = (\mu(f_{k\ell}) - \mu(f_k)\mu(f_\ell) : k, \ell = 1, \dots, d),$$

we have

$$\begin{aligned} \|\text{Cov}(\lambda) - \text{Cov}(\nu)\| &\leq d \max_{k, \ell} (|\lambda(f_{k\ell}) - \nu(f_{k\ell})| + |\lambda(f_k)\lambda(f_\ell) - \nu(f_k)\nu(f_\ell)|) \\ &\leq d \max_{k, \ell} (|\lambda(f_{k\ell}) - \nu(f_{k\ell})| + |\lambda(f_k)||\lambda(f_\ell) - \nu(f_\ell)| \\ &\quad + |\nu(f_\ell)||\lambda(f_k) - \nu(f_k)|) \\ &\leq 3d \text{TV}(\lambda, \nu), \end{aligned}$$

using the fact that  $|\lambda(f_k)| \leq 1$  and  $|\nu(f_k)| \leq 1$  for all  $k$ . □

The following compares the mean and covariance matrix of a distribution before and after transformation by a function.

**Lemma 15.** *Let  $\lambda$  be a Borel distribution with compact support  $A \subset \mathbb{R}^D$  and consider a measurable function  $g : A \mapsto \mathbb{R}^D$  such that  $\|g(x) - x\| \leq \eta$ . Then*

$$\|\mathbb{E}(\lambda^g) - \mathbb{E}(\lambda)\| \leq \eta, \quad \|\text{Cov}(\lambda^g) - \text{Cov}(\lambda)\| \leq \eta(\text{diam}(A) + \eta).$$

*Proof.* Take  $X \sim \lambda$  and let  $Y = g(X) \sim \lambda^g$ . For the means, by Jensen's inequality,

$$\|\mathbb{E}X - \mathbb{E}Y\| \leq \mathbb{E}\|X - Y\| = \mathbb{E}\|X - g(X)\| \leq \eta.$$

For the covariances, we first note that

$$\text{Cov}(X) - \text{Cov}(Y) = \frac{1}{2}(\text{Cov}(X - Y, X + Y) + \text{Cov}(X + Y, X - Y)),$$

where  $\text{Cov}(U, V) := \mathbb{E}((U - \mathbb{E}U)(V - \mathbb{E}V)^T)$  is the cross-covariance of random vectors  $U$  and  $V$ . By Jensen's inequality, the fact  $\|\mathbf{u}\mathbf{v}^T\| = \|\mathbf{u}\|\|\mathbf{v}\|$  for any pair of vectors  $\mathbf{u}, \mathbf{v}$ , and then the Cauchy-Schwarz inequality,

$$\|\text{Cov}(U, V)\| \leq \mathbb{E}(\|U - \mathbb{E}U\| \cdot \|V - \mathbb{E}V\|) \leq \mathbb{E}(\|U - \mathbb{E}U\|^2)^{1/2} \cdot \mathbb{E}(\|V - \mathbb{E}V\|^2)^{1/2}.$$

Hence,

$$\begin{aligned} \|\text{Cov}(X) - \text{Cov}(Y)\| &\leq \|\text{Cov}(X - Y, X + Y)\| \\ &\leq \mathbb{E}[\|X - Y - \mathbb{E}X + \mathbb{E}Y\|^2]^{1/2} \mathbb{E}[\|X + Y - \mathbb{E}X - \mathbb{E}Y\|^2]^{1/2}. \end{aligned}$$

By the fact that the mean minimizes the mean-squared error,

$$\mathbb{E}[\|X - Y - \mathbb{E}X + \mathbb{E}Y\|^2]^{1/2} \leq \mathbb{E}[\|X - Y\|^2]^{1/2} = \mathbb{E}[\|X - g(X)\|^2]^{1/2} \leq \eta.$$

In the same vein, letting  $\mathbf{z} \in \mathbb{R}^D$  be such that  $\|\mathbf{x} - \mathbf{z}\| \leq \frac{1}{2} \text{diam}(A)$  for all  $\mathbf{x} \in A$ , we get

$$\begin{aligned} \mathbb{E}[\|X + Y - \mathbb{E}X - \mathbb{E}Y\|^2]^{1/2} &\leq \mathbb{E}[\|X - \mathbf{z} + Y - \mathbf{z}\|^2]^{1/2} \\ &\leq \mathbb{E}[\|X - \mathbf{z}\|^2]^{1/2} + \mathbb{E}[\|Y - \mathbf{z}\|^2]^{1/2} \\ &\leq 2\mathbb{E}[\|X - \mathbf{z}\|^2]^{1/2} + \mathbb{E}[\|X - g(X)\|^2]^{1/2} \\ &\leq \text{diam}(A) + \eta. \end{aligned}$$

Using the triangle inequality twice. From this, we conclude.  $\square$

Next we compare the covariance matrix of the uniform distribution on a small piece of smooth surface with that of the uniform distribution on the projection of that piece onto a nearby tangent subspace.

**Lemma 16.** *There is a constant  $C_{16} > 0$  depending only on  $d$  such that the following is true. Take  $S \in \mathcal{S}_d(\kappa)$ ,  $r < \frac{1}{C_{16}\kappa}$  and  $\mathbf{x} \in \mathbb{R}^D$  such that  $\text{dist}(\mathbf{x}, S) \leq r$ . Let  $\mathbf{s} = P_S(\mathbf{x})$  and  $T = T_S(\mathbf{s})$ . If  $\boldsymbol{\zeta}$  and  $\boldsymbol{\xi}$  are the means, and  $\mathbf{M}$  and  $\mathbf{N}$  are the covariance matrices, of the uniform distributions on  $S \cap B(\mathbf{x}, r)$  and  $T \cap B(\mathbf{x}, r)$ , respectively, then*

$$\|\boldsymbol{\zeta} - \boldsymbol{\xi}\| \leq C_{16}\kappa r^2, \quad \|\mathbf{M} - \mathbf{N}\| \leq C_{16}\kappa r^3.$$

*Proof.* We focus on proving the bound on the covariances, and leave the bound on the means—whose proof is both similar and simpler—as an exercise to the reader. Let  $T = T_S(\mathbf{s})$ ,  $B_r = B(\mathbf{x}, r)$  and  $g = P_T$  for short. Let  $A = S \cap B_r$  and  $A' = T \cap B_r$ .

First, applying Lemma 15, assuming  $\kappa r$  is sufficiently small, we get that

$$\|\text{Cov}(\lambda_A) - \text{Cov}(\lambda_A^g)\| \leq \frac{\kappa}{2}r^2 \left(2r + \frac{\kappa}{2}r^2\right) \leq 2\kappa r^3, \quad (38)$$

because  $\text{diam}(A) \leq 2r$  and  $\|g(\mathbf{x}) - \mathbf{x}\| \leq \frac{\kappa}{2}\|\mathbf{x} - \mathbf{s}\|^2 \leq \frac{\kappa}{2}r^2$  for all  $\mathbf{x} \in A$  by (10).

Let  $\lambda_{g(A)}$  denote the uniform distribution on  $g(A)$ .  $\lambda_A^g$  and  $\lambda_{g(A)}$  are both supported on  $B_r$ , so that applying Lemma 14 with proper scaling, we get

$$\|\text{Cov}(\lambda_A^g) - \text{Cov}(\lambda_{g(A)})\| \leq 3dr^2 \text{TV}(\lambda_A^g, \lambda_{g(A)}).$$

When  $\kappa r$  is small enough, we know that  $g$  is 1-Lipschitz, and by Lemma 5,  $g^{-1}$  is well-defined and is  $(1 + \kappa r)$ -Lipschitz on  $A'$ . Hence, by Lemma 8 and the fact that  $\dim(A) = d$ , we have

$$\text{TV}(\lambda_A^g, \lambda_{g(A)}) \leq 2((1 + \kappa r)^d - 1) \leq 4d\kappa r,$$

using the inequality  $(1 + x)^d \leq 1 + 2dx$ , valid for any  $x \in [0, 1/d]$  and any  $d \geq 1$ .

Noting that  $\lambda_{A'}$  is also supported on  $B_r$ , applying Lemma 14 with proper scaling, we get

$$\|\text{Cov}(\lambda_{g(A)}) - \text{Cov}(\lambda_{A'})\| \leq 3dr^2 \text{TV}(\lambda_{g(A)}, \lambda_{A'}),$$

with

$$\text{TV}(\lambda_{g(A)}, \lambda_{A'}) \leq 4 \frac{\text{vol}(g(A) \triangle A')}{\text{vol}(A')} \leq 4C_9\kappa(r + \kappa r^2) \leq 8C_9\kappa r,$$

by Lemma 7 and Lemma 9, and assuming that  $\kappa r$  is small enough.

By the triangle inequality,

$$\begin{aligned} \|\mathbf{M} - \mathbf{N}\| &= \|\text{Cov}(\lambda_A) - \text{Cov}(\lambda_{A'})\| \\ &\leq \|\text{Cov}(\lambda_A) - \text{Cov}(\lambda_A^g)\| + \|\text{Cov}(\lambda_A^g) - \text{Cov}(\lambda_{g(A)})\| \\ &\quad + \|\text{Cov}(\lambda_{g(A)}) - \text{Cov}(\lambda_{A'})\| \\ &\leq 2\kappa r^3 + 12d^2\kappa r^3 + 24dC_9\kappa r^3. \end{aligned}$$

From this, we conclude. □

Next is a lemma on the estimation of a covariance matrix. The result is a simple consequence of the matrix Hoeffding inequality of Tropp (2012). Note that simply bounding the operator norm by the Frobenius norm, and then applying the classical Hoeffding inequality (Hoeffding, 1963) would yield a bound sufficient for our purposes, but this is a good opportunity to use a more recent and sophisticated result.

**Lemma 17.** *Let  $\mathbf{C}_m$  denote the empirical covariance matrix based on an i.i.d. sample of size  $m$  from a distribution with support in  $B(0, r) \subset \mathbb{R}^d$  with covariance  $\Sigma$ . Then*

$$\mathbb{P}(\|\mathbf{C}_m - \Sigma\| > r^2 t) \leq 4d \exp\left(-\frac{mt}{16} \min\left(\frac{t}{32}, \frac{m}{d}\right)\right), \quad \forall t \geq 0. \quad (43)$$

*Proof.* Without loss of generality, we assume that the distribution has zero mean and is now supported on  $B(0, 2r)$ . In fact, by a simple rescaling, we may also assume that  $r = 1$ . Let  $\mathbf{x}_1, \dots, \mathbf{x}_m$  denote the sample, with  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$ . We have

$$\mathbf{C}_m = \mathbf{C}_m^* - \frac{1}{m} \bar{\mathbf{x}} \bar{\mathbf{x}}^T,$$

where

$$\mathbf{C}_m^* := \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T, \quad \bar{\mathbf{x}} := \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i.$$

Note that

$$\|\mathbf{C}_m - \Sigma\| \leq \|\mathbf{C}_m^* - \Sigma\| + \frac{1}{m} \|\bar{\mathbf{x}}\|^2.$$



Applying the union bound and then Hoeffding's inequality to each coordinate—which is in  $[-2, 2]$ —we get

$$\mathbb{P}(\|\bar{\mathbf{x}}\| > t) \leq \sum_{j=1}^d \mathbb{P}(|\bar{x}_j| > t/\sqrt{d}) \leq 2d \exp\left(-\frac{mt^2}{8d}\right).$$

Noting that  $\frac{1}{m}(\mathbf{x}_i \mathbf{x}_i^T - \boldsymbol{\Sigma})$ ,  $i = 1, \dots, m$ , are independent, zero-mean, self-adjoint matrices with spectral norm bounded by  $4/m$ , we may apply the matrix Hoeffding inequality (Tropp, 2012, Th. 1.3), to get

$$\mathbb{P}(\|\mathbf{C}_m^* - \boldsymbol{\Sigma}\| > t) \leq 2d \exp\left(-\frac{t^2}{8\sigma^2}\right), \quad \sigma^2 := m(4/m)^2 = 16/m.$$

Applying the union bound and using the previous inequalities, we arrive at

$$\begin{aligned} \mathbb{P}(\|\mathbf{C}_m - \boldsymbol{\Sigma}\| > t) &\leq \mathbb{P}(\|\mathbf{C}_m^* - \boldsymbol{\Sigma}\| > t/2) + \mathbb{P}(\|\bar{\mathbf{x}}\| > \sqrt{mt/2}) \\ &\leq 2d \exp\left(-\frac{mt^2}{512}\right) + 2d \exp\left(-\frac{m^2 t}{16d}\right) \\ &\leq 4d \exp\left(-\frac{mt}{16} \min\left(\frac{t}{32}, \frac{m}{d}\right)\right), \end{aligned}$$

which concludes the proof.  $\square$

#### 6.1.4 PROJECTIONS

We relate below the difference of two orthogonal projections with the largest principal angle between the corresponding subspaces.

**Lemma 18.** *For two affine non-null subspaces  $T, T'$ ,*

$$\|\mathbf{P}_T - \mathbf{P}_{T'}\| = \begin{cases} \sin \theta_{\max}(T, T'), & \text{if } \dim(T) = \dim(T'), \\ 1, & \text{otherwise.} \end{cases}$$

*Proof.* For two affine subspaces  $T, T' \subset \mathbb{R}^D$  of same dimension, let  $\frac{\pi}{2} \geq \theta_1 \geq \dots \geq \theta_D \geq 0$ , denote the principal angles between them. By (Stewart and Sun, 1990, Th. I.5.5), the singular values of  $\mathbf{P}_T - \mathbf{P}_{T'}$  are  $\{\sin \theta_j : j = 1, \dots, q\}$ , so that  $\|\mathbf{P}_T - \mathbf{P}_{T'}\| = \max_j \sin \theta_j = \sin \theta_1 = \sin \theta_{\max}(T, T')$ . Suppose now that  $T$  and  $T'$  are of different dimension, say  $\dim(T) > \dim(T')$ . We have  $\|\mathbf{P}_T - \mathbf{P}_{T'}\| \leq \|\mathbf{P}_T\| \vee \|\mathbf{P}_{T'}\| = 1$ , since  $\mathbf{P}_T$  and  $\mathbf{P}_{T'}$  are orthogonal projections and therefore positive semidefinite with operator norm equal to 1. Let  $L = \mathbf{P}_T(T')$ . Since  $\dim(L) \leq \dim(T') < \dim(T)$ , there is  $\mathbf{u} \in T \cap L^\perp$  with  $\mathbf{u} \neq 0$ . Then  $\mathbf{v}^\top \mathbf{u} = \mathbf{P}_T(\mathbf{v})^\top \mathbf{u} = 0$  for all  $\mathbf{v} \in T'$ , implying that  $\mathbf{P}_{T'}(\mathbf{u}) = 0$  and consequently  $(\mathbf{P}_T - \mathbf{P}_{T'})\mathbf{u} = \mathbf{u}$ , so that  $\|\mathbf{P}_T - \mathbf{P}_{T'}\| \geq 1$ .  $\square$

The lemma below is a perturbation result for eigenspaces and widely known as the  $\sin \Theta$  Theorem of Davis and Kahan (1970). See also (Luxburg, 2007, Th. 7) or (Stewart and Sun, 1990, Th. V.3.6).

**Lemma 19.** *Let  $\mathbf{M}$  be positive semi-definite with eigenvalues  $\beta_1 \geq \beta_2 \geq \dots$ . Suppose that  $\beta_d - \beta_{d+1} > 0$ . Then for any other positive semi-definite matrix  $\mathbf{N}$ ,*

$$\|\mathbf{P}_{\mathbf{N}}^{(d)} - \mathbf{P}_{\mathbf{M}}^{(d)}\| \leq \frac{\sqrt{2}\|\mathbf{N} - \mathbf{M}\|}{\beta_d - \beta_{d+1}},$$

where  $\mathbf{P}_{\mathbf{M}}^{(d)}$  and  $\mathbf{P}_{\mathbf{N}}^{(d)}$  denote the orthogonal projections onto the top  $d$  eigenvectors of  $\mathbf{M}$  and  $\mathbf{N}$ , respectively.

### 6.1.5 INTERSECTIONS

We start with an elementary result on points near the intersection of two affine subspaces.

**Lemma 20.** *Take any two linear subspaces  $T_1, T_2 \subset \mathbb{R}^D$ . For any point  $\mathbf{t}_1 \in T_1 \setminus T_2$ , we have*

$$\text{dist}(\mathbf{t}_1, T_2) \geq \text{dist}(\mathbf{t}_1, T_1 \cap T_2) \sin \theta_{\min}(T_1, T_2).$$

*Proof.* We may reduce the problem to the case where  $T_1 \cap T_2 = \{0\}$ . Indeed, let  $\tilde{T}_1 = T_1 \cap T_2^\perp$ ,  $\tilde{T}_2 = T_1^\perp \cap T_2$  and  $\tilde{\mathbf{t}}_1 = \mathbf{t}_1 - P_{T_1 \cap T_2}(\mathbf{t}_1)$ . Then

$$\|\mathbf{t}_1 - P_{T_2}(\mathbf{t}_1)\| = \|\tilde{\mathbf{t}}_1 - P_{\tilde{T}_2}(\tilde{\mathbf{t}}_1)\|, \quad \|\mathbf{t}_1 - P_{T_1 \cap T_2}(\mathbf{t}_1)\| = \|\tilde{\mathbf{t}}_1\|,$$

and

$$\sin \theta_{\min}(T_1, T_2) = \sin \theta_{\min}(\tilde{T}_1, \tilde{T}_2).$$

So assume that  $T_1 \cap T_2 = \{0\}$ . By (Afriat, 1957, Th. 10.1), the angle formed by  $\mathbf{t}_1$  and  $P_{T_2}(\mathbf{t}_1)$  is at least as large as the smallest principal angle between  $T_1$  and  $T_2$ , which is  $\theta_{\min}(T_1, T_2)$  since  $T_1 \cap T_2 = \{0\}$ . From this the result follows immediately.  $\square$

The following result says that a point cannot be close to two compact and smooth surfaces intersecting at a positive angle without being close to their intersection. Note that the constant there cannot be solely characterized by  $\kappa$ , as it also depends on the separation between the surfaces away from their intersection.

**Lemma 21.** *Suppose  $S_1, S_2 \in \mathcal{S}_d(\kappa)$  intersect at a strictly positive angle and that  $\text{reach}(S_1 \cap S_2) \geq 1/\kappa$ . Then there is a constant  $C_{21}$  such that*

$$\text{dist}(\mathbf{x}, S_1 \cap S_2) \leq C_{21} \max \{ \text{dist}(\mathbf{x}, S_1), \text{dist}(\mathbf{x}, S_2) \}, \quad \forall \mathbf{x} \in \mathbb{R}^D.$$

*Proof.* Assume the result is not true, so there is a sequence  $(\mathbf{x}_n) \subset \mathbb{R}^D$  such that  $\text{dist}(\mathbf{x}_n, S_1 \cap S_2) > n \max_k \text{dist}(\mathbf{x}_n, S_k)$ . Because the surfaces are bounded, we may assume WLOG that the sequence is bounded. Then  $\text{dist}(\mathbf{x}_n, S_1 \cap S_2)$  is bounded, which implies that  $\max_k \text{dist}(\mathbf{x}_n, S_k) = O(1/n)$ . This also forces  $\text{dist}(\mathbf{x}_n, S_1 \cap S_2) \rightarrow 0$ . Indeed, otherwise there is a constant  $C > 0$  and a subsequence  $(\mathbf{x}_{n'}) \subset (\mathbf{x}_n)$  such that  $\text{dist}(\mathbf{x}_{n'}, S_1 \cap S_2) \geq C$ . Since  $(\mathbf{x}_{n'})$  is bounded, there is a subsequence  $(\mathbf{x}_{n''}) \subset (\mathbf{x}_{n'})$  that converges, and by the fact that  $\max_k \text{dist}(\mathbf{x}_{n''}, S_k) = o(1)$ , and by compactness of  $S_k$ , the limit is necessarily in  $S_1 \cap S_2$ , which is incompatible with the fact that  $\text{dist}(\mathbf{x}_{n''}, S_1 \cap S_2) \geq C$ . Since  $\text{dist}(\mathbf{x}_n, S_1 \cap S_2) \rightarrow 1$ , we have

$$n \max_k \text{dist}(\mathbf{x}_n, S_k) < \text{dist}(\mathbf{x}_n, S_1 \cap S_2) = o(1). \quad (47)$$

Assume  $n$  is large enough that  $\text{dist}(\mathbf{x}_n, S_1 \cap S_2) < 1/\kappa$  and let  $\mathbf{s}_n^k$  be the projection of  $\mathbf{x}_n$  onto  $S_k$ , and  $\mathbf{s}_n^\dagger$  the projection of  $\mathbf{x}_n$  onto  $S_1 \cap S_2$ . Let  $T_k = T_{S_k}(\mathbf{s}_n^\dagger)$  and note that  $\theta_{\min}(T_1, T_2) \geq \theta := \theta(S_1, S_2)$ —see definition (6). Let  $\mathbf{t}_n^k$  be the projection of  $\mathbf{s}_n^k$  onto  $T_k$ . Assume WLOG that  $\|\mathbf{t}_n^1 - \mathbf{s}_n^1\| \geq \|\mathbf{t}_n^2 - \mathbf{s}_n^2\|$ . Let  $\mathbf{t}_n$  denote the projection of  $\mathbf{t}_n^1$  onto  $T_1 \cap T_2$ , and then let  $\mathbf{s}_n = P_{S_1 \cap S_2}(\mathbf{t}_n)$ .

By (47), we have

$$n \max_k \|\mathbf{x}_n - \mathbf{s}_n^k\| \leq \|\mathbf{x}_n - \mathbf{s}_n^\dagger\| = o(1). \quad (48)$$

We start with the RHS:

$$\|\mathbf{x}_n - \mathbf{s}_n^\dagger\| = \min_{\mathbf{s} \in S_1 \cap S_2} \|\mathbf{x}_n - \mathbf{s}\| \leq \|\mathbf{x}_n - \mathbf{s}_n\|, \quad (49)$$

and first show that  $\|\mathbf{x}_n - \mathbf{s}_n\| = o(1)$  too. We use the triangle inequality multiple times in what follows. We have

$$\|\mathbf{x}_n - \mathbf{s}_n\| \leq \|\mathbf{x}_n - \mathbf{s}_n^1\| + \|\mathbf{s}_n^1 - \mathbf{t}_n^1\| + \|\mathbf{t}_n^1 - \mathbf{t}_n\| + \|\mathbf{t}_n - \mathbf{s}_n\|. \quad (50)$$

From (48),  $\|\mathbf{x}_n - \mathbf{s}_n^1\| \leq \|\mathbf{x}_n - \mathbf{s}_n^\dagger\| = o(1)$ , so that  $\|\mathbf{s}_n^1 - \mathbf{s}_n^\dagger\| = o(1)$ , and by (10),

$$\|\mathbf{s}_n^1 - \mathbf{t}_n^1\| \leq \frac{\kappa}{2} \|\mathbf{s}_n^1 - \mathbf{s}_n^\dagger\|^2 \leq \kappa(\|\mathbf{s}_n^1 - \mathbf{x}_n\|^2 + \|\mathbf{x}_n - \mathbf{s}_n^\dagger\|^2) = o(1). \quad (51)$$

We also have

$$\|\mathbf{t}_n^1 - \mathbf{t}_n\| = \min_{\mathbf{t} \in T_1 \cap T_2} \|\mathbf{t}_n^1 - \mathbf{t}\| \leq \|\mathbf{t}_n^1 - \mathbf{s}_n^\dagger\| \leq \|\mathbf{t}_n^1 - \mathbf{s}_n^1\| + \|\mathbf{s}_n^1 - \mathbf{x}_n\| + \|\mathbf{x}_n - \mathbf{s}_n^\dagger\| = o(1). \quad (52)$$

Finally,

$$\|\mathbf{t}_n - \mathbf{s}_n\| = \min_{\mathbf{s} \in S_1 \cap S_2} \|\mathbf{t}_n - \mathbf{s}\| \leq \|\mathbf{t}_n - \mathbf{s}_n^\dagger\| \leq \|\mathbf{t}_n - \mathbf{t}_n^1\| + \|\mathbf{t}_n^1 - \mathbf{s}_n^\dagger\| = o(1).$$

We now proceed. The last upper bound is rather crude. Indeed, using (12) for  $S = S_1 \cap S_2$  and  $\mathbf{s} = \mathbf{s}_n^\dagger$ , and noting that  $T_{S_1 \cap S_2}(\mathbf{s}_n^\dagger) = T_1 \cap T_2$ , and using the fact that  $\|\mathbf{t}_n - \mathbf{s}_n^\dagger\| = o(1)$ , we get

$$\|\mathbf{t}_n - \mathbf{s}_n\| \leq \kappa \|\mathbf{t}_n - \mathbf{s}_n^\dagger\|^2 \leq \kappa(\|\mathbf{t}_n - \mathbf{s}_n\| + \|\mathbf{s}_n - \mathbf{x}_n\| + \|\mathbf{x}_n - \mathbf{s}_n^\dagger\|)^2.$$

Using (49),

$$\|\mathbf{t}_n - \mathbf{s}_n\| \leq \kappa(\|\mathbf{t}_n - \mathbf{s}_n\| + 2\|\mathbf{s}_n - \mathbf{x}_n\|)^2 \leq 5\kappa\|\mathbf{x}_n - \mathbf{s}_n\|^2, \quad (53)$$

eventually, since  $\|\mathbf{t}_n - \mathbf{s}_n\| = o(1)$ .

Combining (49), (50), (51) and (53), we get

$$\|\mathbf{x}_n - \mathbf{s}_n\| \leq \|\mathbf{x}_n - \mathbf{s}_n^1\| + O(\|\mathbf{x}_n - \mathbf{s}_n^1\|^2 + \|\mathbf{x}_n - \mathbf{s}_n\|^2) + \|\mathbf{t}_n^1 - \mathbf{t}_n\| + O(\|\mathbf{x}_n - \mathbf{s}_n\|^2),$$

which leads to

$$\|\mathbf{x}_n - \mathbf{s}_n\| \leq 2\|\mathbf{x}_n - \mathbf{s}_n^1\| + 2\|\mathbf{t}_n^1 - \mathbf{t}_n\|, \quad (54)$$

when  $n$  is large enough. Using this bound in (48) combined with (49), we get

$$\|\mathbf{t}_n^1 - \mathbf{t}_n\| \geq \frac{n-2}{2} \max_k \|\mathbf{x}_n - \mathbf{s}_n^k\|.$$

We then have

$$\begin{aligned}
 \max_k \|\mathbf{x}_n - \mathbf{s}_n^k\| &\geq \frac{1}{2} \|\mathbf{s}_n^1 - \mathbf{s}_n^2\| \\
 &\geq \frac{1}{2} (\|\mathbf{t}_n^1 - \mathbf{t}_n^2\| - \|\mathbf{s}_n^1 - \mathbf{t}_n^1\| - \|\mathbf{s}_n^2 - \mathbf{t}_n^2\|) \\
 &\geq \frac{1}{2} \text{dist}(\mathbf{t}_n^1, T_2) - \|\mathbf{s}_n^1 - \mathbf{t}_n^1\|,
 \end{aligned}$$

with

$$\|\mathbf{s}_n^1 - \mathbf{t}_n^1\| = O(\|\mathbf{x}_n - \mathbf{s}_n^1\|^2 + \|\mathbf{x}_n - \mathbf{s}_n^\dagger\|^2) = O(\|\mathbf{x}_n - \mathbf{s}_n\|^2) = O(\|\mathbf{t}_n^1 - \mathbf{t}_n\|^2),$$

due (in the same order) to (51), (48)-(49), and (54). Recalling that  $\|\mathbf{t}_n^1 - \mathbf{t}_n\| = \text{dist}(\mathbf{t}_n^1, T_1 \cap T_2)$ , we conclude that

$$\text{dist}(\mathbf{t}_n^1, T_2) = O(1/n) \text{dist}(\mathbf{t}_n^1, T_1 \cap T_2) + O(1) \text{dist}(\mathbf{t}_n^1, T_1 \cap T_2)^2.$$

However, by Lemma 20, we have  $\text{dist}(\mathbf{t}_n^1, T_2) \geq (\sin \theta) \text{dist}(\mathbf{t}_n^1, T_1 \cap T_2)$ , so that dividing by  $\text{dist}(\mathbf{t}_n^1, T_1 \cap T_2)$  above leads to

$$1 = O(1/n) + O(1) \text{dist}(\mathbf{t}_n^1, T_1 \cap T_2),$$

which is in contradiction with the fact that  $\text{dist}(\mathbf{t}_n^1, T_2) \leq \|\mathbf{t}_n^1 - \mathbf{t}_n\| = o(1)$ , established in (52).  $\square$

### 6.1.6 COVARIANCES NEAR AN INTERSECTION

The following compares a covariance matrix of a distribution supported on the union of two smooth surfaces with that of the projection of that distribution on tangent subspaces.

**Lemma 22.** *Consider  $S_1, S_2 \in \mathcal{S}_d(\kappa)$  intersecting at a positive angle, with  $\text{reach}(S_1 \cap S_2) \geq 1/\kappa$ . Fix  $\mathbf{s} \in S_1$  and let  $A = B(\mathbf{s}, r) \cap (S_1 \cup S_2)$ , where  $r < 1/\kappa$ . Let  $\mathbf{C}$  denote the empirical covariance of an i.i.d. sample from a distribution  $\nu$  supported on  $A$  of size  $m$  and define  $\Sigma = \text{Cov}(\nu)$ . Then*

$$\mathbb{P}(\|\mathbf{C}_m - \Sigma\| > r^2 t + 3\kappa r^3) \leq C_{22} \exp\left(-\frac{mt}{C_{22}} \min(t, m)\right), \quad \forall t \geq 0,$$

where  $C_{22}$  depends only on  $d$ .

*Proof.* Let  $T_1 = T_{S_1}(\mathbf{s})$ , and if  $\text{dist}(\mathbf{s}, S_2) \leq r$ , let  $\mathbf{s}_2 = P_{S_2}(\mathbf{s})$  and  $T_2 = T_{S_2}(\mathbf{s}_2)$ . Define  $g : A \mapsto T_1 \cup T_2$  where  $g(\mathbf{x}) = P_{T_1}(\mathbf{x})$  if  $\mathbf{x} \in S_1$ , and  $g(\mathbf{x}) = P_{T_2}(\mathbf{x})$  otherwise. By (10),  $\|g(\mathbf{x}) - \mathbf{x}\| \leq \frac{\kappa}{2} r^2$  for all  $\mathbf{x} \in A$ . Let  $\Sigma^g = \text{Cov}(\nu^g)$ . Also, let  $\nu_m$  denote the sample distribution and note that  $\mathbf{C}_m = \text{Cov}(\nu_m)$ . Let  $\mathbf{C}_m^g = \text{Cov}(\nu_m^g)$ . Note that  $\nu_m^g$  is the empirical distribution of an i.i.d. sample of size  $m$  from  $\nu^g$ , so the notation is congruent. If  $\text{dist}(\mathbf{s}, S_2) > r$ ,  $\nu^g$  is supported on  $B(\mathbf{s}, r) \cap T_1$ , which is a  $d$ -dimensional ball of radius  $r$ , so that (43) applies to  $\mathbf{C}_m^g$  and  $\Sigma^g$ . If  $\text{dist}(\mathbf{s}, S_2) \leq r$ ,  $\nu^g$  is supported on

$$(B(\mathbf{s}, r) \cap T_1) \cup (B(\mathbf{s}_2, r) \cap T_2) \subset B(\mathbf{s}, 2r) \cap (T_1 \cup T_2) \subset B(\mathbf{s}, 2r) \cap (T_1 + T_2),$$

where the latter is a ball in dimension  $d' := \dim(T_1 + T_2) \leq 2d$  of radius  $2r$ , so that (43)—with  $2r$  in place of  $r$  and  $2d$  in place of  $d$ —applies to  $\mathbf{C}_m^g$  and  $\Sigma^g$ . Then, by the triangle inequality and Lemma 15,

$$\begin{aligned} \|\mathbf{C}_m - \Sigma\| &\leq \|\mathbf{C}_m - \mathbf{C}_m^g\| + \|\mathbf{C}_m^g - \Sigma^g\| + \|\Sigma^g - \Sigma\| \\ &\leq 2\kappa r^3 + \kappa^2 r^4 / 2 + \|\mathbf{C}_m^g - \Sigma^g\|, \end{aligned}$$

as in (38).  $\square$

Here is a continuity result.

**Lemma 23.** *Let  $T_1$  and  $T_2$  be two linear subspaces of same dimension  $d$ . For  $\mathbf{x} \in T_1$ , denote by  $\Sigma_T(\mathbf{x})$  the covariance matrix of the uniform distribution over  $B(\mathbf{x}, r) \cap (T_1 \cup T_2)$ . Then, for all  $\mathbf{x}, \mathbf{y} \in T_1$ ,*

$$\|\Sigma_T(\mathbf{x}) - \Sigma_T(\mathbf{y})\| \leq \begin{cases} 5dr \|\mathbf{x} - \mathbf{y}\|, & \text{if } d \geq 2, \\ r \|\mathbf{x} - \mathbf{y}\| + 2\sqrt{r^3 \|\mathbf{x} - \mathbf{y}\|}, & \text{if } d = 1. \end{cases}$$

We note that, when  $d = 1$ , we also have the bound

$$\|\Sigma_T(\mathbf{x}) - \Sigma_T(\mathbf{y})\| \leq 4\gamma r \|\mathbf{x} - \mathbf{y}\|, \quad \text{if } \min(\text{dist}(\mathbf{x}, T_2), \text{dist}(\mathbf{y}, T_2)) \leq r\sqrt{1 - 1/\gamma^2}.$$

*Proof.* Assume without loss of generality that  $r = 1$ . Let  $A_{\mathbf{x}}^j = B(\mathbf{x}, 1) \cap T_j$  for any  $\mathbf{x}$  and  $j = 1, 2$ . By Lemma 14 and then Lemma 7, we have

$$\begin{aligned} \|\Sigma_T(\mathbf{x}) - \Sigma_T(\mathbf{y})\| &= \|\text{Cov}(\lambda_{A_{\mathbf{x}}^1 \cup A_{\mathbf{x}}^2}) - \text{Cov}(\lambda_{A_{\mathbf{y}}^1 \cup A_{\mathbf{y}}^2})\| \\ &\leq 3d \text{TV}(\lambda_{A_{\mathbf{x}}^1 \cup A_{\mathbf{x}}^2}, \lambda_{A_{\mathbf{y}}^1 \cup A_{\mathbf{y}}^2}) \\ &\leq 12d \frac{\text{vol}((A_{\mathbf{x}}^1 \cup A_{\mathbf{x}}^2) \Delta (A_{\mathbf{y}}^1 \cup A_{\mathbf{y}}^2))}{\text{vol}((A_{\mathbf{x}}^1 \cup A_{\mathbf{x}}^2) \cup (A_{\mathbf{y}}^1 \cup A_{\mathbf{y}}^2))} \\ &\leq 12d \frac{\text{vol}(A_{\mathbf{x}}^1 \Delta A_{\mathbf{y}}^1) + \text{vol}(A_{\mathbf{x}}^2 \Delta A_{\mathbf{y}}^2)}{\text{vol}(A_{\mathbf{x}}^1)}. \end{aligned}$$

Note that  $A_{\mathbf{x}}^2 = B(\mathbf{x}_2, \eta) \cap T_2$  where  $\mathbf{x}_2 := P_{T_2}(\mathbf{x})$  and  $\eta := \sqrt{1 - \|\mathbf{x} - \mathbf{x}_2\|^2} \wedge 1$ . Similarly,  $A_{\mathbf{y}}^2 = B(\mathbf{y}_2, \delta) \cap T_2$  where  $\mathbf{y}_2 := P_{T_2}(\mathbf{y})$  and  $\delta := \sqrt{1 - \|\mathbf{y} - \mathbf{y}_2\|^2} \wedge 1$ . Therefore, applying Lemma 12, we get

$$\frac{\text{vol}(A_{\mathbf{x}}^1 \Delta A_{\mathbf{y}}^1)}{2 \text{vol}(A_{\mathbf{x}}^1)} \leq 1 - (1 - t)_+^d, \quad (64)$$

and assuming WLOG that  $\delta \leq \eta$  (i.e.,  $\|\mathbf{y} - \mathbf{y}_2\| \geq \|\mathbf{x} - \mathbf{x}_2\|$ ) and after proper scaling, we get

$$\frac{\text{vol}(A_{\mathbf{x}}^2 \Delta A_{\mathbf{y}}^2)}{2 \text{vol}(A_{\mathbf{x}}^1)} \leq \zeta := \eta^d - (\eta - t_2)_+^d \wedge \delta^d, \quad (65)$$

where  $t := \|\mathbf{x} - \mathbf{y}\|$  and  $t_2 := \|\mathbf{x}_2 - \mathbf{y}_2\|$ . Note that  $t_2 \leq t$  by the fact that  $P_{T_2}$  is 1-Lipschitz.

For (64), we have  $1 - (1 - t)_+^d \leq dt$ . This is obvious when  $t \geq 1$ , while when  $t \leq 1$  it is obtained using the fact that, for any  $0 \leq a < b \leq 1$ ,

$$b^d - a^d = (b - a)(b^{d-1} + ab^{d-2} + \dots + a^{d-2}b + a^{d-1}) \leq db^{d-1}(b - a) \leq d(b - a). \quad (66)$$

For (65), we consider several cases.

- When  $\eta \leq t_2$ , then  $\zeta = \eta^d \leq \eta \leq t_2 \leq t$ .
- When  $t_2 < \eta \leq t_2 + \delta$ , then  $\zeta = \eta^d - (\eta - t_2)^d \leq dt_2 \leq dt$ , by (66).
- When  $\eta \geq t_2 + \delta$  and  $d \geq 2$ , we have

$$\begin{aligned}
 \zeta &= \eta^d - \delta^d \leq d\eta^{d-1}(\eta - \delta) \leq d\eta(\eta - \delta) \leq d(\eta^2 - \delta^2) \\
 &= d(\|\mathbf{y} - \mathbf{y}_2\|^2 \wedge 1 - \|\mathbf{x} - \mathbf{x}_2\|^2 \wedge 1) \\
 &= d(\|\mathbf{y} - \mathbf{y}_2\| \wedge 1 + \|\mathbf{x} - \mathbf{x}_2\| \wedge 1)(\|\mathbf{y} - \mathbf{y}_2\| \wedge 1 - \|\mathbf{x} - \mathbf{x}_2\| \wedge 1) \\
 &\leq 2d(t + t_2) \leq 4dt,
 \end{aligned}$$

where (66) was applied in the first line and the triangle inequality was applied in the last line, in the form of

$$\|\mathbf{y} - \mathbf{y}_2\| \leq \|\mathbf{y} - \mathbf{x}\| + \|\mathbf{x} - \mathbf{x}_2\| + \|\mathbf{x}_2 - \mathbf{y}_2\| = \|\mathbf{x} - \mathbf{x}_2\| + t + t_2.$$

- When  $\eta \geq t_2 + \delta$  and  $d = 1$ , we have

$$\zeta = \eta - \delta \leq \sqrt{\|\mathbf{y} - \mathbf{y}_2\|^2 \wedge 1 - \|\mathbf{x} - \mathbf{x}_2\|^2 \wedge 1} \leq 2\sqrt{t},$$

using (70), preceded by the fact that, for any  $0 \leq a < b \leq 1$ ,

$$0 \leq \sqrt{1-a} - \sqrt{1-b} = \frac{b-a}{\sqrt{1-a} + \sqrt{1-b}} \leq \frac{b-a}{\sqrt{1-b} + b-a} \leq \frac{b-a}{\sqrt{b-a}} = \sqrt{b-a}.$$

In summary, when  $d \geq 2$ , we can therefore bound  $\|\Sigma_T(\mathbf{x}) - \Sigma_T(\mathbf{y})\|$  by  $dt + 4dt = 5dt$ ; and when  $d = 1$ , we bound that by  $t + 2\sqrt{t}$ .  $\square$

The following is in some sense a converse to Lemma 23, in that we lower-bound the distance between covariance matrices near an intersection of linear subspaces. Note that the covariance matrix does not change when moving parallel to the intersection; however, it does when moving perpendicular to the intersection.

**Lemma 24.** *Let  $T_1$  and  $T_2$  be two linear subspaces of same dimension  $d$  with  $\theta_{\min}(T_1, T_2) \geq \theta_0 > 0$ . Fix a unit norm vector  $\mathbf{v} \in T_1 \cap (T_1 \cap T_2)^\perp$ . With  $\Sigma_T(h\mathbf{v})$  denoting the covariance of the uniform distribution over  $B(h\mathbf{v}, r) \cap (T_1 \cup T_2)$ , we have*

$$\inf_h \sup_\ell \|\Sigma_T(h\mathbf{v}) - \Sigma_T(\ell\mathbf{v})\| \geq r^2/C_{24},$$

where the infimum is over  $0 \leq h \leq 1/\sin \theta_0$  and the supremum over  $\max(0, h - 1/2) \leq \ell \leq \min(1/\sin \theta_0, h + 1/2)$ , and  $C_{24} \geq 1$  depends only on  $d$  and  $\theta_0$ .

*Proof.* Assume without loss of generality that  $r = 1$ . If the statement of the lemma is not true, there are subspaces  $T_1$  and  $T_2$  of same dimension  $d$ , a unit length vector  $\mathbf{v} \in T_1 \cap (T_1 \cap T_2)^\perp$  and  $0 \leq h \leq 1/\sin \theta_0$ , such that

$$\Sigma_T(\ell\mathbf{v}) = \Sigma_T(h\mathbf{v}) \text{ for all } \max(0, h - 1/2) \leq \ell \leq \min(1/\sin \theta_0, h + 1/2). \quad (71)$$

This is because  $\ell \mapsto \Sigma_T(\ell \mathbf{v})$  is continuous once  $T_1, T_2$  and  $\mathbf{v}$  are chosen. By projecting onto  $(T_1 \cap T_2)^\perp$ , we may assume that  $T_1 \cap T_2 = \{0\}$  without loss of generality. Let  $\theta = \angle(\mathbf{v}, T_2)$  and note that  $\theta \geq \theta_0$  since  $T_1 \cap T_2 = \{0\}$ . Define  $\mathbf{u} = (\mathbf{v} - P_{T_2} \mathbf{v}) / \sin \theta$  and also  $\mathbf{w} = P_{T_2} \mathbf{v} / \cos \theta$  when  $\theta < \pi/2$ , and  $\mathbf{w} \in T_2$  is any vector perpendicular to  $\mathbf{v}$  when  $\theta = \pi/2$ .  $B(h\mathbf{v}, 1) \cap T_1$  is the  $d$ -dimensional ball of  $T_1$  of radius 1 and center  $h\mathbf{v}$ , while—using Pythagoras theorem— $B(h\mathbf{v}, 1) \cap T_2$  is the  $d$ -dimensional ball of  $T_2$  of radius  $t := (1 - (h \sin \theta)^2)^{1/2}$  and center  $(h \cos \theta) \mathbf{w}$ . Let  $\lambda$  denote the uniform distribution over  $B(h\mathbf{v}, 1) \cap (T_1 \cup T_2)$  and  $\lambda_k$  the uniform distribution over  $B(h\mathbf{v}, 1) \cap T_k$ . Note that  $\lambda = \alpha \lambda_1 + (1 - \alpha) \lambda_2$ , where

$$\alpha := \frac{\text{vol}(B(h\mathbf{v}, 1) \cap T_1)}{\text{vol}(B(h\mathbf{v}, 1) \cap (T_1 \cup T_2))} = \frac{\text{vol}(B(h\mathbf{v}, 1) \cap T_1)}{\text{vol}(B(h\mathbf{v}, 1) \cap T_1) + \text{vol}(B((h \cos \theta) \mathbf{w}, t) \cap T_2)} = \frac{1}{1 + t^d}.$$

By the law of total covariance,

$$\text{Cov}(\lambda) = \alpha \text{Cov}(\lambda_1) + (1 - \alpha) \text{Cov}(\lambda_2) + \alpha(1 - \alpha)(\mathbb{E}(\lambda_1) - \mathbb{E}(\lambda_2))(\mathbb{E}(\lambda_1) - \mathbb{E}(\lambda_2))^\top. \quad (72)$$

with  $\mathbb{E}(\lambda_1) = h\mathbf{v}$  and  $\mathbb{E}(\lambda_2) = h \cos(\theta) \mathbf{w}$ , and  $\text{Cov}(\lambda_1) = c \mathbf{P}_{T_1}$  and  $\text{Cov}(\lambda_2) = t^2 c \mathbf{P}_{T_2}$ , by Lemma 13. Hence,

$$\Sigma_T(h\mathbf{v}) = \alpha c \mathbf{P}_{T_1} + (1 - \alpha) t^2 c \mathbf{P}_{T_2} + \alpha(1 - \alpha)(1 - t^2) \mathbf{u} \mathbf{u}^\top,$$

using the fact that  $\mathbf{v} - (\cos \theta) \mathbf{w} = (\sin \theta) \mathbf{u}$  and the definition of  $t$ . Let  $\theta_1 = \theta_{\max}(T_1, T_2)$  and let  $\mathbf{v}_1 \in T_1$  be of unit length and such that  $\angle(\mathbf{v}_1, T_2) = \theta_1$ . Then for any  $0 \leq h, \ell \leq 1 / \sin \theta_0$ , we have

$$\|\Sigma_T(h\mathbf{v}) - \Sigma_T(\ell \mathbf{v})\| \geq |\mathbf{v}_1^\top \Sigma_T(h\mathbf{v}) \mathbf{v}_1 - \mathbf{v}_1^\top \Sigma_T(\ell \mathbf{v}) \mathbf{v}_1| = |f(t_h) - f(t_\ell)|, \quad (73)$$

where  $t_h := (1 - (h \sin \theta)^2)^{1/2}$  and

$$f(t) = \frac{c}{1 + t^d} + \frac{ct^{d+2}(\cos \theta_1)^2}{1 + t^d} + \frac{t^d(1 - t^2)(\mathbf{u}^\top \mathbf{v}_1)^2}{(1 + t^d)^2}.$$

It is easy to see that the interval

$$I_h = \{t_\ell : (h - 1/2)_+ \leq \ell \leq (1 / \sin \theta_0) \wedge (h + 1/2)\}$$

is non empty. Because of (71) and (73),  $f(t)$  is constant over  $t \in I_h$ , but this is not possible since  $f$  is a rational function not equal to a constant and therefore cannot be constant over an interval of positive length.  $\square$

We now look at the eigenvalues of the covariance matrix.

**Lemma 25.** *Let  $T_1$  and  $T_2$  be two linear subspaces of same dimension  $d$ . For  $\mathbf{x} \in T_1$ , denote by  $\Sigma_T(\mathbf{x})$  the covariance matrix of the uniform distribution over  $B(\mathbf{x}, 1) \cap (T_1 \cup T_2)$ . Then, for all  $\mathbf{x} \in T_1$ ,*

$$c(1 - (1 - \delta^2(\mathbf{x}))_+^{d/2}) \leq \beta_d(\Sigma_T(\mathbf{x})), \quad \beta_1(\Sigma_T(\mathbf{x})) \leq c + \delta^2(\mathbf{x})(1 - \delta^2(\mathbf{x}))_+^{d/2}, \quad (75)$$

$$\frac{c}{8}(1 - \cos \theta_{\max}(T_1, T_2))^2(1 - \delta^2(\mathbf{x}))_+^{d/2+1} \leq \beta_{d+1}(\Sigma_T(\mathbf{x})) \leq (c + \delta^2(\mathbf{x}))(1 - \delta^2(\mathbf{x}))_+^{d/2}, \quad (76)$$

where  $c := 1/(d + 2)$  and  $\delta(\mathbf{x}) := \text{dist}(\mathbf{x}, T_2)$ .

*Proof.* Applying (72), we have

$$\boldsymbol{\Sigma}_T(\mathbf{x}) = \alpha c \mathbf{P}_{T_1} + (1 - \alpha) c t^2 \mathbf{P}_{T_2} + \alpha(1 - \alpha)(\mathbf{x} - \mathbf{x}_2)(\mathbf{x} - \mathbf{x}_2)^\top, \quad (77)$$

where  $\mathbf{x}_2 := P_{T_2}(\mathbf{x})$  and  $\alpha := (1 + t^d)^{-1}$  with  $t := (1 - \delta^2(\mathbf{x}))_+^{1/2}$ . Because all the matrices in this display are positive semidefinite, we have

$$\beta_d(\boldsymbol{\Sigma}_T(\mathbf{x})) \geq \alpha c \|\mathbf{P}_{T_1}\| = \alpha c,$$

with  $\alpha \geq 1 - t^d$ . And because of the triangle inequality, we have

$$\beta_1(\boldsymbol{\Sigma}_T(\mathbf{x})) \leq \alpha c \|\mathbf{P}_{T_1}\| + (1 - \alpha) c t^2 \|\mathbf{P}_{T_2}\| + \alpha(1 - \alpha) \|\mathbf{x} - \mathbf{x}_2\|^2 \leq c + \alpha(1 - \alpha) \delta^2(\mathbf{x}),$$

with  $\alpha(1 - \alpha) \leq t^d$ . Hence, (75) is proved.

For the upper bound in (76), by Weyl's inequality (Stewart and Sun, 1990, Cor. IV.4.9) and the fact that  $\beta_{d+1}(\mathbf{P}_{T_1}) = 0$ , and then the triangle inequality, we get

$$\begin{aligned} \beta_{d+1}(\boldsymbol{\Sigma}_T(\mathbf{x})) &\leq \|\boldsymbol{\Sigma}_T(\mathbf{x}) - \alpha c \mathbf{P}_{T_1}\| \\ &\leq c(1 - \alpha) t^2 \|\mathbf{P}_{T_2}\| + \alpha(1 - \alpha) \delta^2(\mathbf{x}) \\ &\leq (1 - \alpha)(c + \delta^2(\mathbf{x})), \end{aligned}$$

and we then use the fact that  $1 - \alpha \leq t^d$ .

For the lower bound in (76), let  $\theta_1 \geq \theta_2 \geq \dots \geq \theta_d$  denote the principal angles between  $T_1$  and  $T_2$ . By the definition of principal angles, there are orthonormal bases for  $T_1$  and  $T_2$ , denoted  $\mathbf{v}_1, \dots, \mathbf{v}_d$  and  $\mathbf{w}_1, \dots, \mathbf{w}_d$ , such that  $\mathbf{v}_j^\top \mathbf{w}_k = \mathbb{I}_{j=k} \cdot \cos \theta_j$ . Take  $\mathbf{u} \in \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_d, \mathbf{w}_1)$ , that is, of the form  $\mathbf{u} = a\mathbf{v}_1 + \mathbf{v} + b\mathbf{w}_1$ , with  $\mathbf{v} \in \text{span}(\mathbf{v}_2, \dots, \mathbf{v}_d)$ . Since  $\mathbf{P}_{T_1} = \mathbf{v}_1 \mathbf{v}_1^\top + \dots + \mathbf{v}_d \mathbf{v}_d^\top$  and  $\mathbf{P}_{T_2} = \mathbf{w}_1 \mathbf{w}_1^\top + \dots + \mathbf{w}_d \mathbf{w}_d^\top$ , when  $\|\mathbf{u}\| = 1$ , we have

$$\frac{1}{c} \mathbf{u}^\top \boldsymbol{\Sigma}_T(\mathbf{x}) \mathbf{u} \geq \alpha \mathbf{u}^\top \mathbf{P}_{T_1} \mathbf{u} + (1 - \alpha) t^2 \mathbf{u}^\top \mathbf{P}_{T_2} \mathbf{u}$$

with

$$\mathbf{u}^\top \mathbf{P}_{T_1} \mathbf{u} = a^2 + \|\mathbf{v}\|^2 + 2ab \cos \theta_1 + b^2 \cos^2 \theta_1 = 1 - b^2 \sin^2 \theta_1 \geq 0,$$

and

$$\mathbf{u}^\top \mathbf{P}_{T_2} \mathbf{u} = b^2 + 2ab \cos \theta_1 + a^2 \cos^2 \theta_1 = (a \cos \theta_1 + b)^2.$$

If  $|b| \leq 1/2$ , then  $\mathbf{u}^\top \mathbf{P}_{T_1} \mathbf{u} \geq 3/4$ , implying that  $\frac{1}{c} \mathbf{u}^\top \boldsymbol{\Sigma}_T(\mathbf{x}) \mathbf{u} \geq 3\alpha/4 \geq 3/8$ . Otherwise,

$$\begin{aligned} \mathbf{u}^\top \mathbf{P}_{T_1} \mathbf{u} + \mathbf{u}^\top \mathbf{P}_{T_2} \mathbf{u} &\geq (a + b \cos \theta_1)^2 + (a \cos \theta_1 + b)^2 \\ &\geq (1 - \cos \theta_1)^2 (a^2 + b^2) \geq (1 - \cos \theta_1)^2 / 4, \end{aligned}$$

and using the fact that  $\alpha \geq 1 - \alpha \geq (1 - \alpha) t^2$ , we get

$$\frac{1}{c} \mathbf{u}^\top \boldsymbol{\Sigma}_T(\mathbf{x}) \mathbf{u} \geq (1 - \alpha) t^2 (1 - \cos \theta_1)^2 / 4.$$

This last bound is always worst. Hence, by the Courant-Fischer theorem (Stewart and Sun, 1990, Cor. IV.4.7), we have

$$\beta_{d+1}(\boldsymbol{\Sigma}_T(\mathbf{x})) \geq \frac{c}{4} (1 - \alpha) t^2 (1 - \cos \theta_1)^2,$$

and we conclude with  $1 - \alpha \geq t^d/2$ .  $\square$



Below are two technical results on the covariance matrix of the uniform distribution on the intersection of a ball and the union of two smooth surfaces, near where the surfaces intersect. The first one is in fact a stepping stone to the second. The latter will enable us to apply Lemma 23.

**Lemma 26.** *There is a constant  $C_{26} \geq 3$  such that the following holds. Consider  $S_1, S_2 \in \mathcal{S}_d(\kappa)$ . Fix  $r \leq \frac{1}{C_{26}\kappa}$ , and for  $\mathbf{s} \in S_1$  with  $\text{dist}(\mathbf{s}, S_2) < 1/\kappa$ , let  $\boldsymbol{\Sigma}(\mathbf{s})$  and  $\boldsymbol{\Sigma}_T(\mathbf{s})$  denote the covariance matrices of the uniform distributions over  $B(\mathbf{s}, r) \cap (S_1 \cup S_2)$  and  $B(\mathbf{s}, r) \cap (T_1 \cup T_2)$ , respectively, where  $T_1 := T_{S_1}(\mathbf{s})$  and  $T_2 := T_{S_2}(P_{S_2}(\mathbf{s}))$ . Then*

$$\|\boldsymbol{\Sigma}(\mathbf{s}) - \boldsymbol{\Sigma}_T(\mathbf{s})\| \leq C_{26}\kappa r^3. \quad (83)$$

*Proof.* We note that it is enough to prove the result when  $r$  is small enough. Let  $\delta = \text{dist}(\mathbf{s}, S_2)$  and let  $\mathbf{s}_2 = P_{S_2}(\mathbf{s})$ , so that  $\|\mathbf{s} - \mathbf{s}_2\| = \delta$ . We first treat the case where  $\delta \leq r$ . Let  $B_r$  be short for  $B(\mathbf{s}, r)$  and define  $A_k = B_r \cap S_k$ ,  $\boldsymbol{\mu}_k = \mathbb{E}(\lambda_{A_k})$  and  $\mathbf{D}_k = \text{Cov}(\lambda_{A_k})$ , for  $k = 1, 2$ . Applying (72), we have

$$\boldsymbol{\Sigma}(\mathbf{s}) = \alpha \mathbf{D}_1 + (1 - \alpha) \mathbf{D}_2 + \alpha(1 - \alpha)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top,$$

where

$$\alpha := \frac{\text{vol}(A_1)}{\text{vol}(A_1) + \text{vol}(A_2)}.$$

Recall that  $T_1 = T_{S_1}(\mathbf{s})$  and  $T_2 = T_{S_2}(\mathbf{s}_2)$ , and define  $A'_k = B_r \cap T_k$ , so that  $B_r \cap (T_1 \cup T_2) = A'_1 \cup A'_2$ . Note that  $\mathbb{E}(\lambda_{A'_1}) = \mathbf{s}$  and  $\mathbb{E}(\lambda_{A'_2}) = \mathbf{s}_2$ , and by Lemma 13,  $\mathbf{D}'_1 := \text{Cov}(\lambda_{A'_1}) = cr^2 \mathbf{P}_{T_1}$  and  $\mathbf{D}'_2 := \text{Cov}(\lambda_{A'_2}) = c(r^2 - \delta^2) \mathbf{P}_{T_2}$ , where  $c := 1/(d+2)$ . Applying (72), we have

$$\boldsymbol{\Sigma}_T(\mathbf{s}) = \alpha' \mathbf{D}'_1 + (1 - \alpha') \mathbf{D}'_2 + \alpha'(1 - \alpha')(\mathbf{s} - \mathbf{s}_2)(\mathbf{s} - \mathbf{s}_2)^\top,$$

where

$$\alpha' := \frac{\text{vol}(A'_1)}{\text{vol}(A'_1) + \text{vol}(A'_2)}.$$

By the triangle inequality, we have

$$\|\alpha \mathbf{D}_1 - \alpha' \mathbf{D}'_1\| \leq |\alpha' - \alpha| \|\mathbf{D}'_1\| + \alpha \|\mathbf{D}_1 - \mathbf{D}'_1\| \leq cr^2 |\alpha' - \alpha| + \|\mathbf{D}_1 - \mathbf{D}'_1\|,$$

and similarly,

$$\|(1 - \alpha) \mathbf{D}_2 - (1 - \alpha') \mathbf{D}'_2\| \leq cr^2 |\alpha' - \alpha| + (1 - \alpha) \|\mathbf{D}_2 - \mathbf{D}'_2\|.$$

Assuming that  $\kappa r \leq 1/C_{16}$ , by Lemma 16, we have  $\|\mathbf{D}_1 - \mathbf{D}'_1\| \vee \|\mathbf{D}_2 - \mathbf{D}'_2\| \leq C_{16}\kappa r^3$ . Because  $|\alpha'(1 - \alpha') - \alpha(1 - \alpha)| \leq |\alpha' - \alpha|$  and the identity  $\|\mathbf{a}\mathbf{a}^\top - \mathbf{b}\mathbf{b}^\top\| \leq (\|\mathbf{a}\| + \|\mathbf{b}\|)\|\mathbf{a} - \mathbf{b}\|$ , we also have

$$\begin{aligned} & \|\alpha(1 - \alpha)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top - \alpha'(1 - \alpha')(\mathbf{s} - \mathbf{s}_2)(\mathbf{s} - \mathbf{s}_2)^\top\| \\ & \leq |\alpha' - \alpha| \|\mathbf{s} - \mathbf{s}_2\|^2 + \alpha(1 - \alpha) 3r (\|\boldsymbol{\mu}_1 - \mathbf{s}\| + \|\boldsymbol{\mu}_2 - \mathbf{s}_2\|). \end{aligned}$$

By Lemma 16, we have  $\|\boldsymbol{\mu}_1 - \mathbf{s}\| \vee \|\boldsymbol{\mu}_2 - \mathbf{s}_2\| \leq C_{16}\kappa r^2$ . Assuming that  $\kappa r \leq 1/3$ ,  $P_{T_k}^{-1}$  is well-defined and  $(1 + \kappa r)$ -Lipschitz on  $S_k \cap B_r$ . And being an orthogonal projection,  $P_{T_k}$  is 1-Lipschitz. Hence, applying Lemma 6, we have

$$1 \leq \frac{\text{vol}(A_k)}{\text{vol}(P_{T_k}(A_k))} \leq 1 + \kappa r, \quad k = 1, 2.$$

Then by Lemma 9,

$$\frac{\text{vol}(P_{T_k}(A_k) \triangle A'_k)}{\text{vol}(A'_k)} \leq C_9\kappa(r + \kappa r^2) \leq 2C_9\kappa r, \quad k = 1, 2,$$

when  $\kappa r$  is small enough, implying

$$1 - 2C_9\kappa r \leq \frac{\text{vol}(P_{T_k}(A_k))}{\text{vol}(A'_k)} \leq 1 + 2C_9\kappa r, \quad k = 1, 2.$$

Hence,

$$1 - C\kappa r \leq 1 - 2C_9\kappa r \leq \frac{\text{vol}(A_k)}{\text{vol}(A'_k)} \leq (1 + \kappa r)(1 + 2C_9\kappa r) \leq 1 + C\kappa r, \quad k = 1, 2,$$

for some constant  $C > 0$ . Since for all  $a, b, a', b' > 0$  we have

$$\left| \frac{a}{a+b} - \frac{a'}{a'+b'} \right| \leq \frac{|a-a'| \vee |b-b'|}{(a+b) \vee (a'+b')} \leq |1 - a/a'| \vee |1 - b/b'|,$$

we then get

$$|\alpha - \alpha'| \leq C\kappa r.$$

Applying the triangle inequality to  $\|\boldsymbol{\Sigma}(\mathbf{s}) - \boldsymbol{\Sigma}_T(\mathbf{s})\|$  and using the bounds above, we conclude.

We assumed above that  $\delta \leq r$ . When  $\delta > r$ ,  $A_2 = A'_2 = \emptyset$  and  $\alpha = \alpha' = 1$ . In this case,  $\boldsymbol{\mu}_2, \mathbf{D}_2, D'_2$  are not well-defined, but one can easily check that the above calculations, and the resulting bound, remain valid. In fact, in this case  $\boldsymbol{\Sigma}(\mathbf{s}) = \mathbf{D}_1$  and  $\boldsymbol{\Sigma}_T(\mathbf{s}) = \mathbf{D}'_1$ .  $\square$

**Lemma 27.** *Consider  $S_1, S_2 \in \mathcal{S}_d(\kappa)$  intersecting at a positive angle, with  $\text{reach}(S_1 \cap S_2) \geq 1/\kappa$ . There is a constant  $C_{27} \geq 1$ , depending on  $S_1$  and  $S_2$ , such that the following holds. For  $\mathbf{s} \in S_1$  with  $\text{dist}(\mathbf{s}, S_2) \leq 2r$ , assuming  $r \leq \frac{1}{C_{27}}$ , let  $\boldsymbol{\Sigma}(\mathbf{s})$  and  $\boldsymbol{\Sigma}_T^0(\mathbf{s})$  denote the covariance matrices of the uniform distributions over  $B(\mathbf{s}, r) \cap (S_1 \cup S_2)$  and  $B(P_{T_1^0}(\mathbf{s}), r) \cap (T_1^0 \cup T_2^0)$ , where  $T_k^0 := T_{S_k}(\mathbf{s}^0)$  and  $\mathbf{s}^0 := P_{S_1 \cap S_2}(\mathbf{s})$ . Then*

$$\|\boldsymbol{\Sigma}(\mathbf{s}) - \boldsymbol{\Sigma}_T^0(\mathbf{s})\| \leq C_{27} r^{5/2}.$$

*Proof.* We use the notation of Lemma 26 and its proof. In addition, we let  $C > 0$  denote a constant that depends only on  $S_1$  and  $S_2$  whose value increases with each appearance.

Define  $\mathbf{t}_1 = P_{T_1^0}(\mathbf{s})$  and  $\mathbf{t}_2 = P_{T_2^0}(\mathbf{t}_1)$ . Let  $\delta_0 = \|\mathbf{t}_1 - \mathbf{t}_2\|$ ,  $\delta_1 = \|\mathbf{s} - \mathbf{t}_1\|$ ,  $\delta_2 = \|\mathbf{s}_2 - \mathbf{t}_2\|$ . Define  $A_1^0 = T_1^0 \cap B(\mathbf{t}_1, r)$  and  $A_2^0 = T_2^0 \cap B(\mathbf{t}_2, r)$ . Note that  $A_2^0 = T_2^0 \cap B(\mathbf{t}_2, \sqrt{r^2 - \delta_0^2})$

when  $\delta_0 \leq r$ , and  $A_2^0 = \emptyset$  when  $\delta_0 > r$ . We have  $\|\mathbf{s} - \mathbf{s}^0\| \leq Cr$  by Lemma 21, which then implies  $\delta_1 \leq Cr^2$  by Lemma 2. Assuming that  $\kappa r$  is small enough,

$$\begin{aligned} \delta_2 = \|\mathbf{s}_2 - \mathbf{t}_2\| &\leq \|P_{S_2}(\mathbf{s}) - P_{T_2^0}(\mathbf{s})\| + \|P_{T_2^0}(\mathbf{s}) - P_{T_2^0}(\mathbf{t}_1)\| \\ &\leq C_4\kappa\|\mathbf{s} - \mathbf{s}^0\|^2 + \|\mathbf{s} - \mathbf{t}_1\| \leq Cr^2, \end{aligned}$$

where in the second line we used Lemma 4 and the fact that any metric projection is 1-Lipschitz. Hence,

$$|\delta_0 - \delta| \leq \left| \|\mathbf{t}_1 - \mathbf{t}_2\| - \|\mathbf{s} - \mathbf{s}_2\| \right| \leq \|\mathbf{t}_1 - \mathbf{s}\| + \|\mathbf{t}_2 - \mathbf{s}_2\| = \delta_1 + \delta_2 \leq Cr^2.$$

Note that  $\mathbb{E}(\lambda_{A_1^0}) = \mathbf{t}_1$  and, by Lemma 13 (with  $c$  denoting the constant defined there),  $\mathbf{D}_1^0 := \text{Cov}(\lambda_{A_1^0}) = cr^2\mathbf{P}_{T_1^0}$ . If  $\delta_0 \leq r$ ,  $\mathbb{E}(\lambda_{A_2^0}) = \mathbf{t}_2$  and  $\mathbf{D}_2^0 := \text{Cov}(\lambda_{A_2^0}) = c(r^2 - \delta_0^2)\mathbf{P}_{T_2^0}$ ; if  $\delta_0 > r$ , then  $A_2^0 = \emptyset$ , and we define  $\mathbf{D}_2^0 = 0$  in that case. Applying (72), we have

$$\Sigma_T^0(\mathbf{s}) = \alpha^0\mathbf{D}_1^0 + (1 - \alpha^0)\mathbf{D}_2^0 + \alpha^0(1 - \alpha^0)(\mathbf{t}_1 - \mathbf{t}_2)(\mathbf{t}_1 - \mathbf{t}_2)^\top,$$

where

$$\alpha^0 := \frac{\text{vol}(A_1^0)}{\text{vol}(A_1^0) + \text{vol}(A_2^0)}.$$

Note that  $\alpha^0 = 1$  when  $\delta_0 > r$  (i.e., when  $A_2^0 = \emptyset$ ). We focus on the case where  $\max(\delta, \delta_0) \leq r$ , which is the most complex one. We have

$$\|\mathbf{D}'_1 - \mathbf{D}_1^0\| = \|cr^2\mathbf{P}_{T_1} - cr^2\mathbf{P}_{T_1^0}\| = cr^2\|\mathbf{P}_{T_1} - \mathbf{P}_{T_1^0}\|,$$

and

$$\|\mathbf{D}'_2 - \mathbf{D}_2^0\| = \|c(r^2 - \delta^2)\mathbf{P}_{T_2} - c(r^2 - \delta_0^2)_+\mathbf{P}_{T_2^0}\| \leq cr^2\|\mathbf{P}_{T_2} - \mathbf{P}_{T_2^0}\| + c|\delta^2 - \delta_0^2|,$$

by the triangle inequality and the fact that  $\|\mathbf{P}_T\| \leq 1$  for any affine subspace  $T$ . By Lemma 3 and Lemma 18, we have

$$\|\mathbf{P}_{T_1} - \mathbf{P}_{T_1^0}\| \leq 6\kappa\|\mathbf{s} - \mathbf{s}^0\| \leq Cr, \quad \|\mathbf{P}_{T_2} - \mathbf{P}_{T_2^0}\| \leq 6\kappa\|\mathbf{s}_2 - \mathbf{s}^0\| \leq Cr,$$

since

$$\|\mathbf{s}_2 - \mathbf{s}^0\| \leq \|\mathbf{s}_2 - \mathbf{s}\| + \|\mathbf{s} - \mathbf{s}^0\| \leq 2r + Cr.$$

And since  $|\delta^2 - \delta_0^2| \leq 2r|\delta - \delta_0| \leq Cr^3$ , we have  $\|\mathbf{D}'_k - \mathbf{D}_k^0\| \leq Cr^3$  for  $k = 1, 2$ . Let  $\omega_d$  denote the volume of the  $d$ -dimensional unit ball. Then

$$\text{vol}(A'_1) = \omega_d r^d, \quad \text{vol}(A'_2) = \omega_d (r^2 - \delta^2)^{d/2}, \quad \text{vol}(A_1^0) = \omega_d r^d, \quad \text{vol}(A_2^0) = \omega_d (r^2 - \delta_0^2)^{d/2},$$

so that

$$\begin{aligned} |\alpha' - \alpha^0| &= \left| \frac{1}{1 + (1 - \delta/r)^{d/2}} - \frac{1}{1 + (1 - \delta_0/r)^{d/2}} \right| \\ &\leq |(1 - \delta/r)^{d/2} - (1 - \delta_0/r)^{d/2}|. \end{aligned}$$

Proceeding as when we bounded  $\zeta$  in the proof of Lemma 23, we get

$$|\alpha' - \alpha^0| \leq d\sqrt{|\delta - \delta_0|/r} \leq C\sqrt{r}.$$

Hence, we proved that

$$\|\Sigma_T(\mathbf{s}) - \Sigma_T^0(\mathbf{s})\| \leq Cr^{5/2}.$$

We use this inequality, (83), and the triangle inequality, to get

$$\|\Sigma(\mathbf{s}) - \Sigma_T^0(\mathbf{s})\| \leq \|\Sigma(\mathbf{s}) - \Sigma_T(\mathbf{s})\| + \|\Sigma_T(\mathbf{s}) - \Sigma_T^0(\mathbf{s})\| \leq Cr^{5/2},$$

which is what we needed to prove.  $\square$

## 6.2 Performance Guarantees for Algorithm 2

We deal with the case where there is no noise, that is,  $\tau = 0$  in (7), so that the data points are  $\mathbf{s}_1, \dots, \mathbf{s}_N$  and sampled exactly on  $S_1 \cup S_2$  according to the uniform distribution. We explain how things change when there is noise, meaning  $\tau > 0$ , in Section 6.4.

We start with some notation. Let  $\Xi_i = \{j \neq i : \mathbf{s}_j \in N_r(\mathbf{s}_i)\}$ , with (random) cardinality  $N_i = |\Xi_i|$ . For  $i \in [n]$ , let  $K_i = 1$  if  $\mathbf{s}_i \in S_1$  and  $= 2$  otherwise, and let  $T_i = T_{S_{K_i}}(\mathbf{s}_i)$ , which is the tangent subspace associated with data point  $\mathbf{s}_i$ . Let  $\mathbf{P}_i$  be short for  $\mathbf{P}_{T_i}$ . Let

$$I_\star = \{i : K_j = K_i, \forall j \in \Xi_i\},$$

or equivalently,

$$I_\star^c = \{i : \exists j \text{ s.t. } K_j \neq K_i \text{ and } \|\mathbf{s}_j - \mathbf{s}_i\| \leq r\}.$$

Thus  $I_\star$  indexes the points whose neighborhoods do not contain points from the other cluster.

We will soon define a constant  $C_\bullet$  and will require the following

$$C_\bullet \geq C, \quad r \leq \varepsilon/C, \quad \varepsilon \leq \eta/C, \quad \eta \leq 1/C, \quad (90)$$

for a large enough constant  $C \geq 1$  that depends only on  $S_1$  and  $S_2$ .

### 6.2.1 A CONCENTRATION BOUND FOR LOCAL COVARIANCES

**Objective.** *We derive a concentration bound for the local covariances.*

Let  $\Xi_i = \{j \neq i : \mathbf{s}_j \in N_r(\mathbf{s}_i)\}$ , with (random) cardinality  $N_i = |\Xi_i|$ . When there is no noise,  $\mathbf{C}_i$  is the sample covariance of  $\{\mathbf{s}_j : j \in \Xi_i\}$ . Also, let  $\Sigma(\mathbf{s})$  denote the covariance matrix of the uniform distribution over  $B(\mathbf{s}, r) \cap (S_1 \cup S_2)$ , and note that  $\Sigma_i := \Sigma(\mathbf{s}_i) = \mathbb{E}(\mathbf{C}_i | \mathbf{s}_i)$ . Given  $N_i$ ,  $\{\mathbf{s}_j : j \in \Xi_i\}$  are uniformly distributed on  $B(\mathbf{s}_i, r) \cap (S_1 \cup S_2)$ , and applying Lemma 22, we get that for any  $t > 0$

$$\mathbb{P}(\|\mathbf{C}_i - \Sigma_i\| > r^2t + 3\kappa r^3 \mid \mathbf{s}_i, N_i) \leq C_{22} \exp\left(-\frac{N_i t}{C_{22}} \min(t, N_i)\right). \quad (91)$$

Assume that  $r < 1/(C_{11}\kappa)$  and let  $n_\dagger := nr^d/(2C_{11}/\chi)$  where  $\chi := \frac{\text{vol}_d(S_1) \wedge \text{vol}_d(S_2)}{\text{vol}_d(S_1) + \text{vol}_d(S_2)}$ . Define

$$\Omega_1 = \{\forall i : N_i > n_\dagger\}.$$

Note that, for any  $i$ ,

$$N_i \geq N'_i := \#\{j \neq i : K_i = K_j \text{ and } \mathbf{s}_j \in B(\mathbf{s}_i, r)\},$$

meaning that  $N'_i$  counts the number of neighboring points from the same surface. We have that  $\#\{i : K_i = k\} \sim \text{Bin}(n, \chi_k)$  where  $\chi_k := \frac{\text{vol}_d(S_k)}{\text{vol}_d(S_1) + \text{vol}_d(S_2)}$ , and using a standard concentration inequality for the binomial distribution, there is a numeric constant  $C > 0$  such that

$$\mathbb{P}(\#\{i : K_i = k\} \geq \chi n/2) \geq 1 - \exp(-\chi n/C).$$

(We also used the fact that  $\chi = \chi_1 \wedge \chi_2$ .) With this, we may apply Lemma 11, to get

$$\begin{aligned} \mathbb{P}(\Omega_1^c) &\leq 2 \exp(-\chi n/C) + C_{11} r^{-d} \exp(-(\chi n/2)r^d/C_{11}) \\ &\leq C n \exp(-nr^d/C), \end{aligned}$$

for some other constant  $C > 0$ . (We used the fact that the left-hand side can be made greater than 1 when  $nr^d \leq 1$  by choosing  $C$  large enough, so that we may focus on the case where  $nr^d \geq 1$ .) Define the event

$$\Omega_2 = \left\{ \forall i : \|\mathbf{C}_i - \boldsymbol{\Sigma}_i\| \leq r^2 \eta / C_\bullet + 3\kappa r^3 \right\},$$

where  $C_\bullet$  will be chosen large enough, but fixed, later on. Assume that  $n_{\ddagger} \geq \eta / C_\bullet$ . Note that  $\eta \leq 1$  by assumption, and we may assume  $n_{\ddagger}$  to be larger than any given constant, for if not, the statement in Theorem 1 is void. Based on the union bound and (91), we have

$$\begin{aligned} \mathbb{P}(\Omega_2^c) &\leq \mathbb{P}(\Omega_2^c | \Omega_1) + \mathbb{P}(\Omega_1^c) \\ &\leq n C_{22} \exp(-n_{\ddagger} (\eta / C_\bullet)^2 / C_{22}) + C n \exp(-nr^d/C) \\ &\leq C n \exp(-nr^d \eta^2 / C), \end{aligned}$$

for another constant  $C$  that depends on  $C_\bullet$  and  $d$ .

### 6.2.2 AWAY FROM THE INTERSECTION

**Objective.** *We show that there is an appropriate choice of parameters as in Theorem 1 such that, away from the intersection—when confined to  $I_\star$ —two points belonging to the same surface and within distance  $\varepsilon$  are neighbors, while two points belonging to different surfaces are not neighbors.*

For  $i \in [n]$ , let  $K_i = 1$  if  $\mathbf{s}_i \in S_1$  and  $= 2$  otherwise. Let  $T_i = T_{S_{K_i}}(\mathbf{s}_i)$ , which is the tangent subspace associated with data point  $\mathbf{s}_i$ . Let

$$I_\star = \{i : K_j = K_i, \forall j \in \Xi_i\},$$

or equivalently,

$$I_\star^c = \{i : \exists j \text{ s.t. } K_j \neq K_i \text{ and } \|\mathbf{s}_j - \mathbf{s}_i\| \leq r\}.$$

By definition,  $I_\star$  indexes the points whose neighborhoods do not contain points from the other cluster.

For  $i \in I_\star$ , let  $\Sigma_{T,i}$  denote the covariance of the uniform distribution on  $T_i \cap B(\mathbf{s}_i, r)$ . Applying Lemma 16, this leads to

$$\|\Sigma_i - \Sigma_{T,i}\| \leq C_{16}\kappa r^3, \quad \forall i \in I_\star. \quad (95)$$

Under  $\Omega_2$ , by the triangle inequality and (95), we have

$$\|\mathbf{C}_i - \Sigma_{T,i}\| \leq \|\mathbf{C}_i - \Sigma_i\| + \|\Sigma_i - \Sigma_{T,i}\| \leq \zeta r^2, \quad \forall i \in I_\star, \quad (96)$$

where

$$\zeta := \eta/C_\bullet + (3 + C_{16})\kappa r. \quad (97)$$

The inequality (96) leads, via the triangle inequality, to the decisive bound

$$\|\mathbf{C}_i - \mathbf{C}_j\| \leq \|\Sigma_{T,i} - \Sigma_{T,j}\| + 2\zeta r^2, \quad \forall i, j \in I_\star. \quad (98)$$

Take  $i, j \in I_\star$  such that  $K_i = K_j$  and  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq \varepsilon$ . Then by Lemma 13, Lemma 18 and Lemma 3, and the triangle inequality, we have

$$\frac{1}{cr^2} \|\Sigma_{T,i} - \Sigma_{T,j}\| = \sin \theta_{\max}(T_i, T_j) \leq 6\kappa \|\mathbf{s}_i - \mathbf{s}_j\| \leq 6\kappa\varepsilon, \quad (99)$$

where  $c := 1/(d+2)$ . This implies that

$$\frac{1}{r^2} \|\mathbf{C}_i - \mathbf{C}_j\| \leq 6c\kappa\varepsilon + 2\zeta,$$

by the triangle inequality and (98). Therefore, if  $\eta > 6c\kappa\varepsilon + 2\zeta$ , then any pair of points indexed by  $i, j \in I_\star$  from the same cluster and within distance  $\varepsilon$  are direct neighbors in the graph built by Algorithm 2.

Take  $i, j \in I_\star$  such that  $K_i \neq K_j$  and  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq \varepsilon$ . By Lemma 21,

$$\max [\text{dist}(\mathbf{s}_i, S_1 \cap S_2), \text{dist}(\mathbf{s}_j, S_1 \cap S_2)] \leq C_{21} \|\mathbf{s}_i - \mathbf{s}_j\|.$$

Let  $\mathbf{z}$  be the mid-point of  $\mathbf{s}_i$  and  $\mathbf{s}_j$ . By convexity and the display above,

$$\text{dist}(\mathbf{z}, S_1 \cap S_2) \leq \frac{1}{2} \text{dist}(\mathbf{s}_i, S_1 \cap S_2) + \frac{1}{2} \text{dist}(\mathbf{s}_j, S_1 \cap S_2) \leq C_{21}\varepsilon.$$

Assuming  $C_{21}\varepsilon < 1/\kappa$ , let  $\mathbf{s} = P_{S_1 \cap S_2}(\mathbf{z})$ . Then, by the triangle inequality again,

$$\max [\|\mathbf{s} - \mathbf{s}_i\|, \|\mathbf{s} - \mathbf{s}_j\|] \leq \text{dist}(\mathbf{z}, S_1 \cap S_2) + \frac{1}{2} \|\mathbf{s}_i - \mathbf{s}_j\| \leq C_{21}\varepsilon + \frac{1}{2}\varepsilon.$$

Let  $T'_i$  denote the tangent subspace of  $S_{K_i}$  at  $\mathbf{s}$  and let  $\Sigma'_i$  be the covariance of the uniform distribution over  $T'_i \cap B(\mathbf{s}, r)$ . Define  $T'_j$  and  $\Sigma'_j$  similarly. Then, as in (99) we have

$$\frac{1}{cr^2} \|\Sigma_{T,i} - \Sigma'_i\| \leq 6\kappa \|\mathbf{s}_i - \mathbf{s}\| \leq 6\kappa(C_{21} + 1/2)\varepsilon,$$

and similarly,

$$\frac{1}{cr^2} \|\Sigma_{T,j} - \Sigma'_j\| \leq 6\kappa(C_{21} + 1/2)\varepsilon.$$

Moreover, by Lemma 13 and Lemma 18,

$$\frac{1}{cr^2} \|\Sigma'_i - \Sigma'_j\| = \sin \theta_{\max}(T'_i, T'_j) \geq \sin \theta_s,$$

where  $\theta_s$  is short for  $\theta(S_1, S_2)$ , defined in (6). Hence, by the triangle inequality,

$$\frac{1}{cr^2} \|\Sigma_{T,i} - \Sigma_{T,j}\| \geq \sin \theta_s - 12c\kappa(C_{21} + 1/2)\varepsilon,$$

and then

$$\frac{1}{r^2} \|\mathbf{C}_i - \mathbf{C}_j\| \geq c \sin \theta_s - 12c\kappa(C_{21} + 1/2)\varepsilon - 2\zeta, \quad (100)$$

by the triangle inequality and (98). Therefore, if  $\eta < c \sin \theta_s - 12c\kappa(C_{21} + 1/2)\varepsilon - 2\zeta$ , then any pair of points indexed by  $i, j \in I_\star$  from different clusters are *not* direct neighbors in the graph built by Algorithm 2.

In summary, we would like to choose  $\eta$  such that

$$c\kappa\varepsilon + 2\zeta < \eta < c \sin \theta_s - 12c\kappa(C_{21} + 1/2)\varepsilon - 2\zeta.$$

Recalling the definition of  $\zeta$  in (97), this holds under (90).

### 6.2.3 THE DIFFERENT CLUSTERS ARE DISCONNECTED IN THE GRAPH

**Objective.** *We show that Step 2 in Algorithm 2 eliminates all points  $i \notin I_\star$ , implying by our conclusion in Section 6.2.2 that the two clusters are not connected in the graph.*

Hence, take  $i \notin I_\star$  with  $K_i = 1$  (say), so that  $\text{dist}(\mathbf{s}_i, S_2) \leq r$ . By Lemma 21, we have  $\text{dist}(\mathbf{s}_i, S_1 \cap S_2) \leq C_{21}r$ . Assume  $r$  is small enough that  $C_{21}\kappa r < 1/2$ . Let  $\mathbf{s}^0 = P_{S_1 \cap S_2}(\mathbf{s}_i)$  and define  $T_k^0 = T_{S_k}(\mathbf{s}^0)$  for  $k \in \{1, 2\}$ . For  $\mathbf{s} \in S_1$  such that  $\text{dist}(\mathbf{s}, S_2) \leq 2r$ , define  $\Sigma_T^0(\mathbf{s})$  as in Lemma 27.

Assuming that  $\mathbf{s}_i \neq \mathbf{s}^0$  (which is true with probability one) and  $\mathbf{s}^0 = 0$  (by translation invariance), let  $h = \|\mathbf{s}_i - \mathbf{s}^0\|$  and  $\mathbf{v} = (\mathbf{s}_i - \mathbf{s}^0)/h$ . Note that  $\mathbf{s}_i = h\mathbf{v}$ . Because  $\mathbf{v} \perp T_1^0 \cap T_2^0$  by Lemma 1, and that  $\theta_{\min}(T_1^0, T_2^0) \geq \theta_s$ , we apply Lemma 24 to find  $\ell \in h \pm r/2$  such that  $\|\Sigma_T^0(\ell\mathbf{v}) - \Sigma_T^0(h\mathbf{v})\| \geq r^2/C_{24}$ , where  $C_{24} \geq 1$  depends only on  $\theta_s$  and  $d$ . Letting  $\tilde{\mathbf{s}} = \ell\mathbf{v}$ , we have  $\|\tilde{\mathbf{s}} - \mathbf{s}_i\| = |h - \ell| \leq r/2$ , so that

$$\text{dist}(\tilde{\mathbf{s}}, S_2) \leq \|\tilde{\mathbf{s}} - \mathbf{s}_i\| + \text{dist}(\mathbf{s}_i, S_2) \leq r/2 + r \leq 3r/2.$$

By Lemma 21, we have  $\text{dist}(\tilde{\mathbf{s}}, S_1 \cap S_2) \leq C_{21}2r < 1/\kappa$ , and consequently,  $P_{S_1 \cap S_2}(\tilde{\mathbf{s}}) = \mathbf{s}^0$ , by Lemma 1. Hence, by the triangle inequality and Lemma 27,

$$\begin{aligned} \|\Sigma(\mathbf{s}_i) - \Sigma(\tilde{\mathbf{s}})\| &\geq \|\Sigma_T^0(\mathbf{s}) - \Sigma_T^0(\tilde{\mathbf{s}})\| - \|\Sigma(\mathbf{s}_i) - \Sigma_T^0(\mathbf{s}_i)\| - \|\Sigma(\tilde{\mathbf{s}}) - \Sigma_T^0(\tilde{\mathbf{s}})\| \\ &\geq r^2/C_{24} - 2C_{27}r^{5/2}. \end{aligned}$$

Take any  $\bar{\mathbf{s}} \in S_1$  such that  $\|\bar{\mathbf{s}} - \tilde{\mathbf{s}}\| \leq r/2$ . By Lemma 23,

$$\|\Sigma_T^0(\bar{\mathbf{s}}) - \Sigma_T^0(\tilde{\mathbf{s}})\| \leq 5dr\|\bar{\mathbf{s}} - \tilde{\mathbf{s}}\| + 2\sqrt{r^3\|\bar{\mathbf{s}} - \tilde{\mathbf{s}}\|} \leq (5d + 2)r^{3/2}\sqrt{\|\bar{\mathbf{s}} - \tilde{\mathbf{s}}\|}.$$

Noting that

$$\text{dist}(\bar{\mathbf{s}}, S_2) \leq \|\bar{\mathbf{s}} - \tilde{\mathbf{s}}\| + \text{dist}(\tilde{\mathbf{s}}, S_2) \leq r/2 + 3r/2 \leq 2r,$$

we apply the triangle inequality and Lemma 27, and we get

$$\begin{aligned} \|\Sigma(\bar{\mathbf{s}}) - \Sigma(\tilde{\mathbf{s}})\| &\leq \|\Sigma_T^0(\bar{\mathbf{s}}) - \Sigma_T^0(\tilde{\mathbf{s}})\| + \|\Sigma(\bar{\mathbf{s}}) - \Sigma_T^0(\bar{\mathbf{s}})\| + \|\Sigma(\tilde{\mathbf{s}}) - \Sigma_T^0(\tilde{\mathbf{s}})\| \\ &\leq (5d+2)r^{3/2}\sqrt{\|\bar{\mathbf{s}} - \tilde{\mathbf{s}}\|} + 2C_{27}r^{5/2}. \end{aligned}$$

When  $\|\bar{\mathbf{s}} - \tilde{\mathbf{s}}\| \leq r/C_{\dagger}$ , we have

$$\begin{aligned} \|\Sigma(\bar{\mathbf{s}}) - \Sigma(\mathbf{s}_i)\| &\geq \|\Sigma(\mathbf{s}_i) - \Sigma(\tilde{\mathbf{s}})\| - \|\Sigma(\bar{\mathbf{s}}) - \Sigma(\tilde{\mathbf{s}})\| \\ &\geq r^2(1/C_{24} - (5d+2)/\sqrt{C_{\dagger}} - 4C_{27}r^{1/2}) \geq r^2/(2C_{24}), \end{aligned}$$

assuming  $C_{\dagger}$  is large enough and  $r$  is small enough. Fixing such a constant  $C_{\dagger} \geq 2$ , define the event

$$\Omega_3 = \bigcap_{k=1}^2 \left\{ \forall \mathbf{s} \in S_k : \#\{i : K_i = k \text{ and } \mathbf{s}_i \in B(\mathbf{s}, r/C_{\dagger})\} \geq 2 \right\}.$$

By the same arguments bounding  $\mathbb{P}(\Omega_1^c)$  in Section 6.2.1, we have that there is a constant  $C$  such that, when  $r \leq 1/(C\kappa)$ ,

$$\mathbb{P}(\Omega_3^c) \leq Cn \exp(-nr^d/C).$$

Under  $\Omega_3$ , there is  $\mathbf{s}_j \in S_1 \cap B(\tilde{\mathbf{s}}, r/C_{\dagger})$ . Taking  $\bar{\mathbf{s}} = \mathbf{s}_j$  above, we have that  $\|\Sigma_j - \Sigma_i\| \geq r^2/(2C_{24})$ . Under  $\Omega_2$ , by the triangle inequality,

$$\begin{aligned} \|\mathbf{C}_j - \mathbf{C}_i\| &\geq \|\Sigma_j - \Sigma_i\| - \|\mathbf{C}_i - \Sigma_i\| - \|\mathbf{C}_j - \Sigma_j\| \\ &\geq r^2/(2C_{24}) - 2r^2\eta/C_{\bullet} \geq r^2/(3C_{24}), \end{aligned}$$

choosing  $C_{\bullet}$  sufficiently large. (Recall that  $\eta \leq 1$  by assumption.) Therefore, choosing  $\eta$  such that  $\eta < 1/(3C_{24})$ , we see that  $\|\mathbf{C}_j - \mathbf{C}_i\| > \eta r^2$ , even though

$$\|\mathbf{s}_j - \mathbf{s}_i\| \leq \|\mathbf{s}_j - \tilde{\mathbf{s}}\| + \|\tilde{\mathbf{s}} - \mathbf{s}_i\| \leq r/C_{\dagger} + r/2 \leq r.$$

#### 6.2.4 THE POINTS THAT ARE REMOVED ARE CLOSE TO THE INTERSECTION

**Objective.** *We show that the points removed by Step 2 are within distance  $Cr$  of the intersection.*

Let  $I_o = \{i \in I_{\star} : \Xi_i \subset I_{\star}\}$ . By our choice of parameters in (90), we see that  $i \in I_o$  is neighbor with any  $j \in \Xi_i$ , so that  $i$  survives Step 2. Hence, the nodes removed at Step 2 are in  $I_o^c = \{i \in I_{\star} : \Xi_i \cap I_{\star}^c \neq \emptyset\} \cup I_{\star}^c$ , with the possibility that some nodes in  $I_o^c$  survive. Now, for any  $i \in I_{\star}^c$ , there is  $j$  with  $K_j \neq K_i$  such that  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq r$ , so by Lemma 21,

$$\text{dist}(\mathbf{s}_i, S_1 \cap S_2) \leq C_{21}\|\mathbf{s}_i - \mathbf{s}_j\| \leq C_{21}r.$$

And for any  $i \in I_o \setminus I_{\star}^c$ , there is  $j \in I_{\star}^c$  such that  $\|\mathbf{s}_j - \mathbf{s}_i\| \leq r$ , by definition of  $\Xi_i$ , so that  $\text{dist}(\mathbf{s}_i, S_1 \cap S_2) \leq C_{21}r + r$ , by the triangle inequality. In fact, we just proved that the last inequality holds for all  $i \in I_o^c$ . Hence, the points removed in Step 2 are within distance  $(C_{21} + 1)r$  of  $S_1 \cap S_2$ .



## 6.2.5 EACH CLUSTER IS (ESSENTIALLY) CONNECTED IN THE GRAPH

**Objective.** *We show that the points that survive Step 2 and belong to the same cluster are connected in the graph, except for possible spurious points within distance  $Cr$  of the intersection.*

Take, for example, the cluster generated from sampling  $S_1$ . Before we start, we recall that  $I_*^c$  was eliminated in Step 2, so that by our choice of  $\eta$  in (90), to show that two points  $\mathbf{s}_i, \mathbf{s}_j$  sampled from  $S_1$  are neighbors it suffices to show that  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq \varepsilon$ .

Define  $R = S_1 \setminus B(S_1 \cap S_2, r')$ , where  $r' := (C_{21} + 2)r$ . By our conclusion in Section 6.2.4, we have that

$$\text{any } \mathbf{s}_i \in S_1 \text{ such that } \text{dist}(\mathbf{s}_i, R) < r \text{ survives Step 2,} \quad (109)$$

so that it suffices to show that the subset of nodes  $\{i : \mathbf{s}_i \in R\}$  is connected.

Define  $W = S_1 \setminus S_2$  and let  $\{W_m\}$  denote the connected components of  $W$ . (Note that there might only be one, for example, when  $S_1$  and  $S_2$  are circles in dimension  $D = 3$ , not in the same plane, that intersect at a single point.) Define  $R_m = W_m \setminus B(S_1 \cap S_2, r')$ . Since  $W = \sqcup_m W_m = S_1 \setminus (S_1 \cap S_2)$ , we have  $R = \sqcup_m R_m$ . In fact, when  $r$  (and therefore  $r'$ ) is small enough,  $R_m$  is connected, and assuming this is the case,  $\{R_m\}$  are the connected components of  $R$ .

Suppose  $R_1$  and  $R_2$  are two adjacent connected components of  $R$ , meaning that  $R_1 \cup R_2 \cup S_2$  is connected. We show that there is at least one pair  $j_1, j_2$  of direct neighbors in the graph such that  $\mathbf{s}_{j_m} \in R_m$ . Take  $\mathbf{s}$  on the connected component of  $S_1 \cap S_2$  separating  $R_1$  and  $R_2$ . We construct two points,  $\mathbf{s}^1 \in R_1$  and  $\mathbf{s}^2 \in R_2$ , that are within distance of order  $\varepsilon$  from  $\mathbf{s}$ , and then select  $\mathbf{s}_{j_m}$  as the closest data point to  $\mathbf{s}^m$ . Let  $T^k = T_{S_k}(\mathbf{s})$  and let  $H$  be the affine subspace parallel to  $(T^1 \cap T^2)^\perp$  passing through  $\mathbf{s}$ . Take  $\mathbf{t}^m \in P_{T^1}(R_m) \cap H \cap \partial B(\mathbf{s}, \varepsilon_1)$ , where  $\varepsilon_1 := \varepsilon/4$ , which exists when  $\varepsilon$  is sufficiently small. Note that  $\mathbf{t}^1, \mathbf{t}^2 \in T^1$  and assuming that  $\varepsilon$  is sufficiently small that  $\varepsilon_1 < 1/\kappa$ , we have  $P_{S_1 \cap S_2}(\mathbf{t}^m) = \mathbf{s}$  by Lemma 1. Define  $\mathbf{s}^m = P_{T^1}^{-1}(\mathbf{t}^m)$  and note that  $\mathbf{s}^1, \mathbf{s}^2 \in S_1$ . Lemma 5 not only justifies this construction when  $\kappa\varepsilon_1 < 1/C_5$ , it also says that  $P_{T^1}^{-1}$  has Lipschitz constant bounded by  $1 + \kappa\varepsilon_1$ , which implies that

$$\|\mathbf{s}^m - \mathbf{s}\| \leq (1 + \kappa\varepsilon_1)\|\mathbf{t}^m - \mathbf{s}\| = (1 + \kappa\varepsilon_1)\varepsilon_1 \leq \varepsilon/3.$$

We also have

$$\begin{aligned} \text{dist}(\mathbf{s}^m, S_1 \cap S_2) &\geq \text{dist}(\mathbf{t}^m, S_1 \cap S_2) - \|\mathbf{s}^m - \mathbf{t}^m\| \\ &= \|\mathbf{t}^m - \mathbf{s}\| - \|\mathbf{s}^m - \mathbf{t}^m\| \\ &\geq \varepsilon_1 - \frac{\kappa}{2}\|\mathbf{s}^m - \mathbf{s}\|^2 \\ &\geq \left(1 - \frac{\kappa}{2}(1 + \kappa\varepsilon_1)^2\varepsilon_1\right)\varepsilon_1 \\ &\geq \varepsilon/5, \end{aligned}$$

when  $\varepsilon$  is sufficiently small. We used (10) in the second inequality. We assume  $r/\varepsilon$  is sufficiently small that  $\varepsilon/5 \geq r' + r$ . Then under  $\Omega_3$ , there are  $j_1, j_2$  such that  $\mathbf{s}_{j_m} \in B(\mathbf{s}^m, r) \cap S_1$ , and by the triangle inequality, we then have that  $\text{dist}(\mathbf{s}_{j_m}, S_1 \cap S_2) \geq \varepsilon/5 - r \geq r'$ , so

that  $\mathbf{s}_{j_m} \in R_m$ . Moreover,

$$\begin{aligned} \|\mathbf{s}_{j_1} - \mathbf{s}_{j_2}\| &\leq \|\mathbf{s}_{j_1} - \mathbf{s}^1\| + \|\mathbf{s}^1 - \mathbf{s}\| + \|\mathbf{s} - \mathbf{s}^2\| + \|\mathbf{s}^2 - \mathbf{s}_{j_2}\| \\ &\leq r + \varepsilon/3 + \varepsilon/3 + r \\ &= \frac{2}{3}\varepsilon + 2r \leq \varepsilon, \end{aligned}$$

assuming  $r/\varepsilon$  is sufficiently small.

Now, we show that the points sampled from any connected component, say  $R_1$ , form a connected subgraph. Take  $\mathbf{s}^1, \dots, \mathbf{s}^M$  an  $r$ -packing of  $R_1$ , so that

$$\bigsqcup_m (R_1 \cap B(\mathbf{s}^m, r/2)) \subset R_1 \subset \bigcup_m (R_1 \cap B(\mathbf{s}^m, r)).$$

Because  $R_1$  is connected,  $\cup_m B(\mathbf{s}^m, r)$  is necessarily connected. Under  $\Omega_1$ , and  $C_{26} \geq 2$ , all the balls  $B(\mathbf{s}^m, r)$ ,  $m = 1, \dots, M$ , contain at least one  $\mathbf{s}_i \in S_1$ , and any such point survives Step 2, because of  $\text{dist}(\mathbf{s}_i, R_1) < r$  and (109). Assume  $r/\varepsilon \leq 1/4$ . Then two points  $\mathbf{s}_i$  and  $\mathbf{s}_j$  such that  $\mathbf{s}_i, \mathbf{s}_j \in B(\mathbf{s}^m, r)$  are connected, since  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq 2r \leq \varepsilon$ . And when  $B(\mathbf{s}^{m_1}, r) \cap B(\mathbf{s}^{m_2}, r) \neq \emptyset$ ,  $\mathbf{s}_i \in B(\mathbf{s}^{m_1}, r)$  and  $\mathbf{s}_j \in B(\mathbf{s}^{m_2}, r)$  are such that  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq 4r \leq \varepsilon$ . Hence, the points sampled from  $R_1$  are connected in the graph under  $\Omega_1$ .

### 6.2.6 SUMMARY AND CONCLUSION

We worked under the assumption that (90), with a constant  $C$  large enough and depending only on  $S_1$  and  $S_2$ . This is compatible with the statement of Theorem 1. In Section 6.2.1, we derived a concentration inequality for local covariances. This lead to a uniform bound (defining the event  $\Omega_2$ ) which holds with probability at least  $1 - Cn \exp(-nr^d \eta^2/C)$  for some constant  $C > 0$ . In Section 6.2.2, we showed that among points away from the intersection (indexed by  $I_\star$ ), those that are in the same cluster are neighbors in the graph if they are within distance  $\varepsilon$ , while those in different clusters cannot be neighbors in the graph. In Section 6.2.3, we showed that Step 2 removes points not indexed by  $I_\star$ , thus implying that the two clusters are not connected in the graph. In the process, we assumed that an event (denoted  $\Omega_3$ ) held, which happens with probability at least  $1 - Cn \exp(-nr^d/C)$  for some constant  $C > 0$ . In Section 6.2.4, we showed that the points removed in Step 2 are within distance  $O(r)$  from the intersection of the surfaces. In Section 6.2.5, we showed that, except for points within distance  $O(r)$  of the intersection, two points in the same cluster are connected in the graph.

We conclude that, after Step 4, Algorithm 2 returns two groups, and each group covers an entire cluster except for points within distance  $O(r)$  of the intersection. In Step 5, each point removed in Step 2 is included in the group which contains its closest point (in Euclidean distance). Theorem 1 is silent on the accuracy of doing this step. So the proof of Theorem 1, as it relates to Algorithm 2, is complete.

### 6.3 Performance Guarantees for Algorithm 3

We keep the same notation as in Section 6.2 and go a little faster here as the arguments are parallel. Let  $d_i$  denote the estimated dimensionality at point  $\mathbf{s}_i$ , meaning the number

of eigenvalues of  $\mathbf{C}_i$  exceeding  $\sqrt{\eta} \|\mathbf{C}_i\|$ . Recall that  $\mathbf{Q}_i$  denotes the orthogonal projection onto the top  $d_i$  eigenvectors of  $\mathbf{C}_i$ .

### 6.3.1 ESTIMATION OF THE INTRINSIC DIMENSION

**Objective.** *We show that, for  $i \in I_\star$ ,  $d_i = d$  (the correct intrinsic dimension) with high probability.*

Take  $i \in I_\star$ . Under  $\Omega_2$ , (96) holds, and applying Weyl's inequality (Stewart and Sun, 1990, Cor. IV.4.9), we have

$$|\beta_m(\mathbf{C}_i) - \beta_m(\boldsymbol{\Sigma}_{T,i})| \leq \zeta r^2, \quad \forall m = 1, \dots, D.$$

By Lemma 13,  $\boldsymbol{\Sigma}_{T,i} = cr^2 \mathbf{P}_i$ , so that  $\beta_m(\boldsymbol{\Sigma}_{T,i}) = cr^2$  when  $m \leq d$  and  $\beta_m(\boldsymbol{\Sigma}_{T,i}) = 0$  when  $m > d$ . Hence,

$$\beta_1(\mathbf{C}_i) \leq (c + \zeta)r^2, \quad \beta_d(\mathbf{C}_i) \geq (c - \zeta)r^2, \quad \beta_{d+1}(\mathbf{C}_i) \leq \zeta r^2.$$

This implies that

$$\frac{\beta_d(\mathbf{C}_i)}{\beta_1(\mathbf{C}_i)} \geq \frac{c - \zeta}{c + \zeta} > \sqrt{\eta}, \quad \frac{\beta_{d+1}(\mathbf{C}_i)}{\beta_1(\mathbf{C}_i)} \leq \frac{\zeta}{c + \zeta} < \sqrt{\eta},$$

when  $\zeta \leq \eta/2$  and  $\eta$  is sufficiently small, which is the case under (90). When this is so,  $d_i = d$  by definition of  $d_i$ . (Note that  $\|\mathbf{C}_i\| = \beta_1(\mathbf{C}_i)$ .)

### 6.3.2 EACH CLUSTER IS (ESSENTIALLY) CONNECTED IN THE GRAPH

**Objective.** *We show that each cluster is connected in the graph, except possibly for some points near the intersection.*

Note that the top  $d$  eigenvectors of  $\boldsymbol{\Sigma}_{T,i}$  generate  $T_i$ . Hence, applying Lemma 19, and (96) again, we get that (recall that  $c = 1/(d+1)$ )

$$\|\mathbf{Q}_i - \mathbf{P}_i\| \leq \frac{\sqrt{2} \zeta r^2}{cr^2} = \zeta' := \sqrt{2}(d+2)\zeta, \quad \forall i \in I_\star.$$

This is the equivalent of (96), which leads to the equivalent of (98):

$$\|\mathbf{Q}_i - \mathbf{Q}_j\| \leq \frac{1}{cr^2} \|\boldsymbol{\Sigma}_{T,i} - \boldsymbol{\Sigma}_{T,j}\| + 2\zeta', \quad \forall i, j \in I_\star,$$

using the fact that  $\boldsymbol{\Sigma}_{T,i} = cr^2 \mathbf{P}_i$ . When  $i, j \in I_\star$  are such that  $K_i = K_j$ , based on (99), we have

$$\|\mathbf{Q}_i - \mathbf{Q}_j\| \leq 6\kappa\varepsilon + 2\zeta'.$$

Hence, when  $\eta > 6\kappa\varepsilon + 2\zeta'$  (which is the case under (90)), two nodes  $i, j \in I_\star$  such that  $K_i = K_j$  and  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq \varepsilon$  are neighbors in the graph. The arguments provided in Section 6.2.5 now apply in exactly the same way to show that nodes  $i \in I_o$  such that  $K_i = 1$  are connected in the graph. The same is true of nodes  $i \in I_o$  such that  $K_i = 2$ .

## 6.3.3 THE CLUSTERS ARE (ESSENTIALLY) DISCONNECTED IN THE GRAPH

**Objective.** *We show that the two sets of nodes  $\{i \in I_o : K_i = 1\}$  and  $\{i \in I_o : K_i = 2\}$ —shown to be connected above—are disconnected in the graph.*

When we take  $i, j \in I_*$  such that  $K_i \neq K_j$ , we have the equivalent of (100), meaning,

$$\|\mathbf{Q}_i - \mathbf{Q}_j\| \geq \sin \theta_s - 12\kappa(C_{21} + 1/2)\varepsilon - 2\zeta'.$$

We choose  $\eta$  smaller than the RHS, so that these nodes are not neighbors in the graph.

We next prove that a point indexed by  $I_*$  is not neighbor to a point near the intersection because of different estimates for the local dimension. Let  $C$  denote a positive constant which increases with each appearance. Take  $\mathbf{s} \in S_1$  such that  $\delta(\mathbf{s}) := \text{dist}(\mathbf{s}, S_2) \leq r$ . Reinststate the notation used in Section 6.2.3, in particular  $\mathbf{s}^0$ ,  $T_1^0$  and  $T_2^0$ , as well as  $\Sigma_T^0(\mathbf{s})$ . Define  $\Sigma_T(\mathbf{s})$  as in Lemma 26. Also, let  $T_1' = T_{S_1}(\mathbf{s})$  and  $T_2' = T_{S_2}(\mathbf{s}^2)$  where  $\mathbf{s}^2 := P_{S_2}(\mathbf{s})$ . By Lemma 26 and Weyl's inequality, we have

$$\beta_{d+1}(\Sigma(\mathbf{s})) \geq \beta_{d+1}(\Sigma_T(\mathbf{s})) - C_{26}r^3, \quad \beta_1(\Sigma(\mathbf{s})) \leq \beta_1(\Sigma_T(\mathbf{s})) + C_{26}r^3,$$

which together with Lemma 25 (and proper scaling), implies that

$$\frac{\beta_{d+1}(\Sigma(\mathbf{s}))}{\beta_1(\Sigma(\mathbf{s}))} \geq \frac{\frac{\varepsilon}{8}(1 - \cos \theta_{\max}(T_1', T_2'))^2(1 - \delta(\mathbf{s})^2/r^2)^{d/2+1} - C_{26}r^3}{c + (\delta(\mathbf{s})/r)^2(1 - \delta(\mathbf{s})^2/r^2)^{d/2} + C_{26}r^3}.$$

Then, by the triangle inequality,

$$\theta_{\max}(T_1', T_2') \geq \theta_{\max}(T_1^0, T_2^0) - \theta_{\max}(T_1', T_1^0) - \theta_{\max}(T_2', T_2^0).$$

By definition,  $\theta_{\max}(T_1^0, T_2^0) \geq \theta_s$ , and by Lemma 3,

$$\theta_{\max}(T_1', T_1^0) \leq \sin^{-1}(6\kappa\|\mathbf{s} - \mathbf{s}^0\| \wedge 1) \leq Cr,$$

and similarly,

$$\theta_{\max}(T_2', T_2^0) \leq \sin^{-1}(6\kappa\|\mathbf{s}^2 - \mathbf{s}^0\| \wedge 1) \leq Cr,$$

because  $\|\mathbf{s} - \mathbf{s}^0\| \leq Cr$  by Lemma 21 and then  $\|\mathbf{s}^2 - \mathbf{s}^0\| \leq r + \|\mathbf{s} - \mathbf{s}^0\| \leq Cr$ . Hence, for  $r$  small enough,  $\theta_{\max}(T_1', T_2') \geq \theta_s/2$ , and furthermore,

$$\frac{\beta_{d+1}(\Sigma(\mathbf{s}))}{\beta_1(\Sigma(\mathbf{s}))} \geq \sqrt{\eta} \quad \text{when} \quad 1 - \delta(\mathbf{s})^2/r^2 \geq \xi := C_*\eta^{1/(d+2)}, \quad (118)$$

where  $C_*$  is a large enough constant, assuming  $r/\eta$  and  $\eta$  are small enough. The same is true for points on  $\mathbf{s} \in S_2$  if we redefine  $\delta(\mathbf{s}) = \text{dist}(\mathbf{s}, S_1)$ . Hence, for  $\mathbf{s}_i$  close enough to the intersection that  $\delta(\mathbf{s}_i)$  satisfies (118),  $d_i > d$ . Then, by Lemma 18,  $\|\mathbf{Q}_i - \mathbf{Q}_j\| = 1$  for any  $j \in I_*$ . By our choice of  $\eta < 1$ , this means that  $i$  and  $j$  are not neighbors.

So the only way  $\{i \in I_* : K_i = 1\}$  and  $\{i \in I_* : K_i = 2\}$  are connected in the graph is if there are  $\mathbf{s}_i \in S_1$  and  $\mathbf{s}_j \in S_2$  such that  $\|\mathbf{s}_i - \mathbf{s}_j\| \leq \varepsilon$  and both  $\delta(\mathbf{s}_i)$  and  $\delta(\mathbf{s}_j)$  fail to satisfy (118). We now show this is not possible. Let  $\Sigma_{T,i} = \Sigma_T(\mathbf{s}_i)$  and define  $\Sigma_{T,j}$  similarly. (These are well-defined when  $\varepsilon < 1/\kappa$ , which we assume.) By Lemma 26, we have

$$\|\Sigma_i - \Sigma_{T,i}\| \leq C_{26}r^3. \quad (119)$$

Define  $\alpha_i = (1 + t_i^d)^{-1}$  with  $t_i := (1 - \delta^2(\mathbf{s}_i))_+^{1/2}$ . As in (77), with the triangle inequality, we get

$$\begin{aligned} \|\boldsymbol{\Sigma}_{T,i} - \alpha_i c r^2 \mathbf{P}_i\| &\leq c(1 - \alpha_i) t_i^2 r^2 + \alpha_i (1 - \alpha_i) \delta^2(\mathbf{s}_i) \\ &\leq 2(1 - \alpha_i) r^2 \leq 2(1 - \delta(\mathbf{s}_i)^2 / r^2)_+^{d/2} r^2 \leq 2\xi^{d/2} r^2, \end{aligned}$$

where the second inequality we used the fact that  $\alpha_i = 1$  if  $\delta(\mathbf{s}_i) > r$ , and in very last inequality comes from  $\delta(\mathbf{s}_i)$  not satisfying (118). Hence, under  $\Omega_2$ , with (119) and the triangle inequality, we get

$$\|\mathbf{C}_i - \alpha_i c r^2 \mathbf{P}_i\| \leq r^2 \eta / C_\bullet + 3\kappa r^3 + 2\xi^{d/2} r^2 + C_{26} r^3.$$

Since  $\beta_d(\mathbf{P}_i) = 1$  and  $\beta_{d+1}(\mathbf{P}_i) = 0$ , by Lemma 19, we have

$$\|\mathbf{Q}_i - \mathbf{P}_i\| \leq \frac{1}{\alpha_i c r^2} [r^2 \eta / C_\bullet + 2\xi^{d/2} r^2 + (3\kappa + C_{26}) r^3] \leq C(\eta / C_\bullet + \xi^{d/2} + r).$$

Similarly,

$$\|\mathbf{Q}_j - \mathbf{P}_j\| \leq C(\eta / C_\bullet + \xi^{d/2} + r).$$

By Lemma 18,  $\|\mathbf{P}_i - \mathbf{P}_j\| = \sin \theta_{\max}(T_i, T_j)$ . Let  $\mathbf{s}^0 = P_{S_1 \cap S_2}(\mathbf{s}_i)$ , and define  $T_1^0$  and  $T_2^0$  as before. Note that  $\|\mathbf{s}_i - \mathbf{s}^0\| \leq C\varepsilon$  by Lemma 21 and the fact that  $\text{dist}(\mathbf{s}_i, S_2) \leq \|\mathbf{s}_i - \mathbf{s}_j\| \leq \varepsilon$ . Also,  $\|\mathbf{s}_j - \mathbf{s}^0\| \leq \|\mathbf{s}_i - \mathbf{s}^0\| + \|\mathbf{s}_j - \mathbf{s}_i\| \leq C\varepsilon$ . We then have

$$\theta_{\max}(T_i, T_j) \geq \theta_{\max}(T_1^0, T_2^0) - \theta_{\max}(T_i, T_1^0) - \theta_{\max}(T_j, T_2^0) \geq \theta_s - C\varepsilon,$$

by Lemma 3. Hence, by the triangle inequality,

$$\begin{aligned} \|\mathbf{Q}_i - \mathbf{Q}_j\| &\geq \|\mathbf{P}_i - \mathbf{P}_j\| - \|\mathbf{Q}_i - \mathbf{P}_i\| - \|\mathbf{Q}_j - \mathbf{P}_j\| \\ &\geq \sin(\theta_s - C\varepsilon) - C(\eta / C_\bullet + \xi^{d/2} + r) \geq \frac{1}{2} \sin \theta_s > \eta, \end{aligned}$$

when the constant in (90) is large enough. (Recall the definition of  $\xi$  above.) Therefore  $i$  and  $j$  are not neighbors, as we needed to show.

#### 6.4 Noisy Case

So far we only dealt with the case where  $\tau = 0$  in (7). When  $\tau > 0$ , a sample point  $\mathbf{x}_i$  is in general different than its corresponding point  $\mathbf{s}_i$  sampled from one of the surfaces. However, when  $\tau/r$  is small, this does not change things much. For one thing, the points are close to each other, since we have  $\|\mathbf{x}_i - \mathbf{s}_i\| \leq \tau$  by assumption, and  $\tau$  is small compared to  $r$ . And the corresponding covariance matrices are also close to each other. To see this, redefine  $\Xi_i = \{j \neq i : \mathbf{x}_j \in N_r(\mathbf{x}_i)\}$  and  $\mathbf{C}_i$  as the sample covariance of  $\{\mathbf{x}_j : j \in \Xi_i\}$ . Let  $\mathbf{D}_i$  denote the sample covariance of  $\{\mathbf{s}_j : j \in \Xi_i\}$ . Let  $X$  be uniform over  $\{\mathbf{x}_j : j \in \Xi_i\}$  and define  $Y = \sum_j \mathbf{s}_j \mathbb{1}_{\{X=\mathbf{x}_j\}}$ . Starting from (33), we have

$$\begin{aligned} \|\mathbf{D}_i - \mathbf{C}_i\| &= \|\text{Cov}(X) - \text{Cov}(Y)\| \\ &\leq \mathbb{E} [\|X - Y\|^2]^{1/2} \cdot \left( \mathbb{E} [\|X - \mathbf{x}_i\|^2]^{1/2} + \mathbb{E} [\|Y - \mathbf{x}_i\|^2]^{1/2} \right) \\ &\leq \tau \cdot (r + (r + \tau)) = r^2(2\tau/r + (\tau/r)^2), \end{aligned}$$

which is small compared to  $r^2$ , and  $r^2$  is the operating scale for covariance matrices in our setting.

Using these facts, the arguments are virtually the same, except for some additional terms due to triangle inequalities, for example,  $\|\mathbf{s}_i - \mathbf{s}_j\| - 2\tau \leq \|\mathbf{x}_i - \mathbf{x}_j\| \leq \|\mathbf{s}_i - \mathbf{s}_j\| + 2\tau$ . In particular, this results in  $\zeta$  in (97) being now redefined as  $\zeta = \frac{3\tau}{r} + \eta/C_\bullet + (3 + C_{16})\kappa r$ . We omit further technical details.

## Acknowledgments

This work was partially supported by grants from the National Science Foundation (DMS 0915160, 0915064, 0956072, 1418386, 1513465). We would like to thank Jan Rataj for helpful discussion around Lemma 3 and Xu Wang for his sharp proofreading. We also gratefully acknowledge the comments, suggestions, and scrutiny of an anonymous referee. We would also like to acknowledge support from the Institute for Mathematics and its Applications (IMA). For one thing, the authors first learned about the research of Goldberg et al. (2009) there, at the *Multi-Manifold Data Modeling and Applications* workshop in the Fall of 2008, and this was the main inspiration for our paper. Also, part of our work was performed while TZ was a postdoctoral fellow at the IMA, and also while EAC and GL were visiting the IMA.

## References

- S. N. Afriat. Orthogonal and oblique projectors and the characteristics of pairs of vector spaces. *Proceedings of the Cambridge Philosophical Society*, 53:800–816, 1957.
- S. Agarwal, J. Lim, L. Zelnik-Manor, P. Perona, D. Kriegman, and S. Belongie. Beyond pairwise clustering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 838–845, 2005.
- S. Agarwal, K. Branson, and S. Belongie. Higher order learning with graphs. In *International Conference on Machine Learning (ICML)*, pages 17–24, 2006.
- E. Arias-Castro. Clustering based on pairwise distances when the data is of mixed dimensions. *IEEE Transactions on Information Theory*, 57(3):1692–1706, 2011.
- E. Arias-Castro, G. Chen, and G. Lerman. Spectral clustering based on local linear approximations. *Electronic Journal of Statistics*, 5:1537–1587, 2011.
- R. Basri and D. W. Jacobs. Lambertian reflectance and linear subspaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(2):218–233, 2003.
- J.-D. Boissonnat, R. Dyer, and A. Ghosh. Constructing intrinsic delaunay triangulations of submanifolds. *arXiv preprint arXiv:1303.6493*, 2013.
- P. Bradley, K. Bennett, and A. Demiriz. Constrained  $k$ -means clustering. Technical Report MSR-TR-2000-65, Microsoft Research, 2000.

- M. R. Brito, E. L. Chávez, A. J. Quiroz, and J. E. Yukich. Connectivity of the mutual k-nearest-neighbor graph in clustering and outlier detection. *Statistics & Probability Letters*, 35(1):33–42, 1997.
- G. Chen and G. Lerman. Spectral curvature clustering (scc). *International Journal of Computer Vision*, 81(3):317–330, 2009a.
- G. Chen and G. Lerman. Foundations of a multi-way spectral clustering framework for hybrid linear modeling. *Foundations of Computational Mathematics*, 9(5):517–558, 2009b.
- G. Chen, S. Atev, and G. Lerman. Kernel spectral curvature clustering (KSCC). In *IEEE International Conference on Computer Vision (ICCV) Workshops*, pages 765–772, 2009.
- A. Cuevas, R. Fraiman, and B. Pateiro-López. On statistical properties of sets fulfilling rolling-type conditions. *Advances in Applied Probability*, 44(2):311–329, 2012.
- C. Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. III. *SIAM Journal on Numerical Analysis*, 7:1–46, 1970.
- E. Elhamifar and R. Vidal. Sparse manifold clustering and embedding. In *Advances in Neural Information Processing Systems (NIPS)*, volume 24, pages 55–63, 2011.
- R. Epstein, P. Hallinan, and A. Yuille.  $5 \pm 2$  eigenimages suffice: An empirical investigation of low-dimensional lighting models. In *IEEE Workshop on Physics-based Modeling in Computer Vision*, pages 108–116, 1995.
- Z. Fan, J. Zhou, and Y. Wu. Multibody grouping by inference of multiple subspaces from high-dimensional data using oriented-frames. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(1):91–105, 2006.
- H. Federer. Curvature measures. *Transactions of the American Mathematical Society*, 93:418–491, 1959.
- Z. Fu, W. Hu, and T. Tan. Similarity based vehicle trajectory clustering and anomaly detection. In *IEEE International Conference on Image Processing*, volume 2, pages 602–605, 2005.
- A. Gionis, A. Hinneburg, S. Papadimitriou, and P. Tsaparas. Dimension induced clustering. In *International Conference on Knowledge Discovery in Data Mining (KDD)*, pages 51–60, 2005.
- A. Goh and R. Vidal. Segmenting motions of different types by unsupervised manifold clustering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–6, 2007.
- A.B. Goldberg, X. Zhu, A. Singh, Z. Xu, and R. Nowak. Multi-manifold semi-supervised learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 169–176, 2009.
- D. Gong, X. Zhao, and G. Medioni. Robust multiple manifolds structure learning. In *International Conference on Machine Learning (ICML)*, pages 321–328, 2012.

- Q. Guo, H. Li, W. Chen, I-F. Shen, and J. Parkkinen. Manifold clustering via energy minimization. In *International Conference on Machine Learning and Applications (ICMLA)*, pages 375–380, 2007.
- G. Haro, G. Randall, and G. Sapiro. Stratification learning: Detecting mixed density and dimensionality in high dimensional point clouds. *Advances in Neural Information Processing Systems (NIPS)*, 19:553, 2007.
- R. Heckel and H. Bölcskei. Noisy subspace clustering via thresholding. In *IEEE International Symposium on Information Theory (ISIT)*, pages 1382–1386, 2013.
- J. Ho, M.-H. Yang, J. Lim, K.-C. Lee, and D. Kriegman. Clustering appearances of objects under varying illumination conditions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 1–11, 2003.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- D. N. Kaslovsky and F. G. Meyer. Optimal tangent plane recovery from noisy manifold samples. arXiv preprint arXiv:1111.4601v2, 2011.
- D. Kushnir, M. Galun, and A. Brandt. Fast multiscale clustering and manifold identification. *Pattern Recognition*, 39(10):1876–1891, 2006.
- K.-C. Lee, J. Ho, and D. J. Kriegman. Acquiring linear subspaces for face recognition under variable lighting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(5):684–698, 2005.
- A. V. Little, Y. M. Jung, and M. Maggioni. Multiscale estimation of intrinsic dimensionality of data sets. In *Manifold Learning and its Applications*, pages 26–33, 2009.
- U. Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- Y. Ma, A. Y. Yang, H. Derksen, and R. Fossum. Estimation of subspace arrangements with applications in modeling and segmenting mixed data. *SIAM Review*, 50(3):413–458, 2008.
- M. Maier, M. Hein, and U. von Luxburg. Optimal construction of  $k$ -nearest-neighbor graphs for identifying noisy clusters. *Theoretical Computer Science*, 410(19):1749–1764, 2009.
- V. J. Martínez and E. Saar. *Statistics of the Galaxy Distribution*. CRC press, Boca Raton, 2002.
- A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. *Advances in Neural Information Processing Systems (NIPS)*, 2:849–856, 2002.
- P. Niyogi, S. Smale, and S. Weinberger. Finding the homology of submanifolds with high confidence from random samples. *Discrete & Computational Geometry*, 39(1):419–441, 2008.



- M. Polito and P. Perona. Grouping and dimensionality reduction by locally linear embedding. *Advances in Neural Information Processing Systems (NIPS)*, 14:1255–1262, 2001.
- L. K. Saul and S. T. Roweis. Think globally, fit locally: unsupervised learning of low dimensional manifolds. *Journal of Machine Learning Research*, 4:119–155, 2003.
- A. Shashua, R. Zass, and T. Hazan. Multi-way clustering using super-symmetric non-negative tensor factorization. In *European Conference on Computer Vision (ECCV)*, pages 595–608, 2006.
- M. Soltanolkotabi and E. J. Candès. A geometric analysis of subspace clustering with outliers. *Annals of Statistics*, 40(4):2195–2238, 2012.
- M. Soltanolkotabi, E. Elhamifar, and E. J. Candès. Robust subspace clustering. *Annals of Statistics*, 42(2):669–699, 2014.
- R. Souvenir and R. Pless. Manifold clustering. In *IEEE International Conference on Computer Vision (ICCV)*, volume 1, pages 648–653, 2005.
- G. W. Stewart and J. G. Sun. *Matrix Perturbation Theory*. Computer Science and Scientific Computing. Academic Press Inc., Boston, MA, 1990.
- M. Tipping and C. Bishop. Mixtures of probabilistic principal component analysers. *Neural Computation*, 11(2):443–482, 1999.
- J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12:389–434, 2012.
- M. C. Tsakiris and R. Vidal. Filtrated spectral algebraic subspace clustering. In *IEEE International Conference on Computer Vision Workshops*, pages 28–36, 2015.
- R. Valdarnini. Detection of non-random patterns in cosmological gravitational clustering. *Astronomy & Astrophysics*, 366:376–386, 2001.
- R. Vidal and Y. Ma. A unified algebraic approach to 2-D and 3-D motion segmentation and estimation. *Journal of Mathematical Imaging and Vision*, 25(3):403–421, 2006.
- G. Walther. Granulometric smoothing. *Annals of Statistics*, 25(6):2273–2299, 1997.
- Y. Wang, Y. Jiang, Y. Wu, and Z.-H. Zhou. Spectral clustering on multiple manifolds. *IEEE Transactions on Neural Networks*, 22(7):1149–1161, 2011.
- Y.-X. Wang and H. Xu. Noisy sparse subspace clustering. In *International Conference on Machine Learning (ICML)*, pages 89–97, 2013.
- L. Zelnik-Manor and P. Perona. Self-tuning spectral clustering. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1601–1608, 2005.
- T. Zhang, A. Szlam, Y. Wang, and G. Lerman. Hybrid linear modeling via local best-fit flats. *International Journal of Computer Vision*, pages 1–24, 2012.