

Fitting a manifold to noisy data

Charles Fefferman, Sergei Ivanov, Yaroslav Kurylev, Matti Lassas,
Hariharan Narayanan

Manifold Hypothesis

- Often, high dimensional data can be modeled to lie in the vicinity of a low dimensional manifold.
- In such cases, the manifold hypothesis is said to hold true.

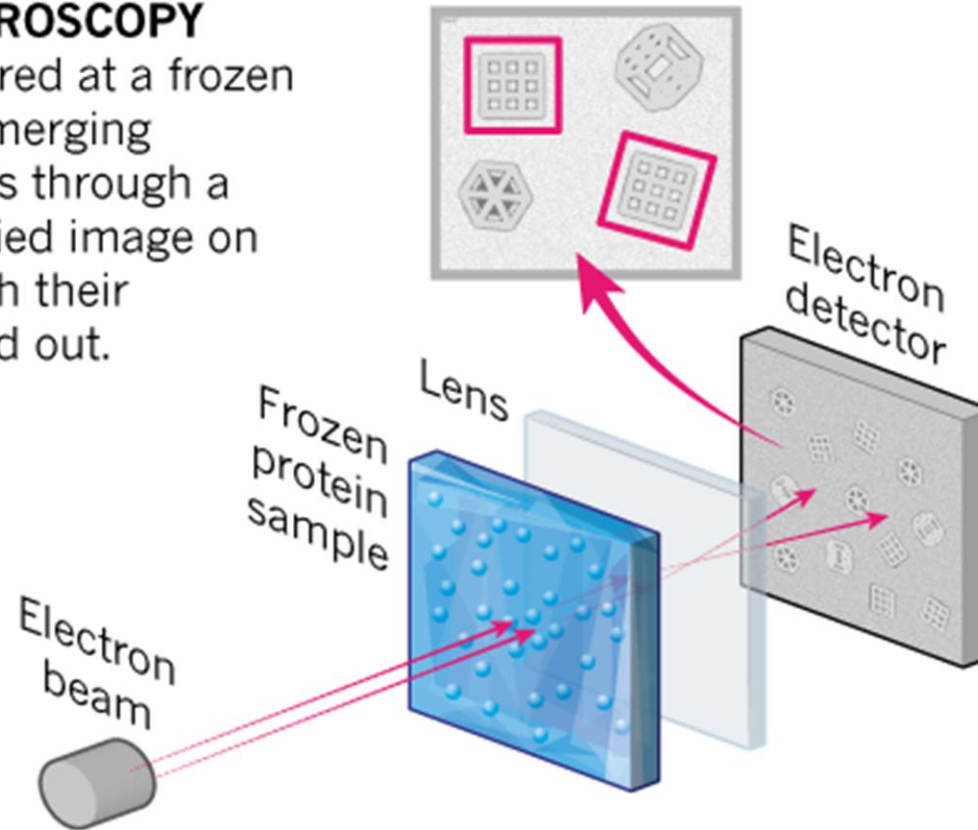
Outline of talk

- An Instance: Cryo EM
- Testing the manifold hypothesis
- Finding a putative manifold.

Measurements can lie near a manifold

CRYO-ELECTRON MICROSCOPY

A beam of electron is fired at a frozen protein solution. The emerging scattered electrons pass through a lens to create a magnified image on the detector, from which their structure can be worked out.

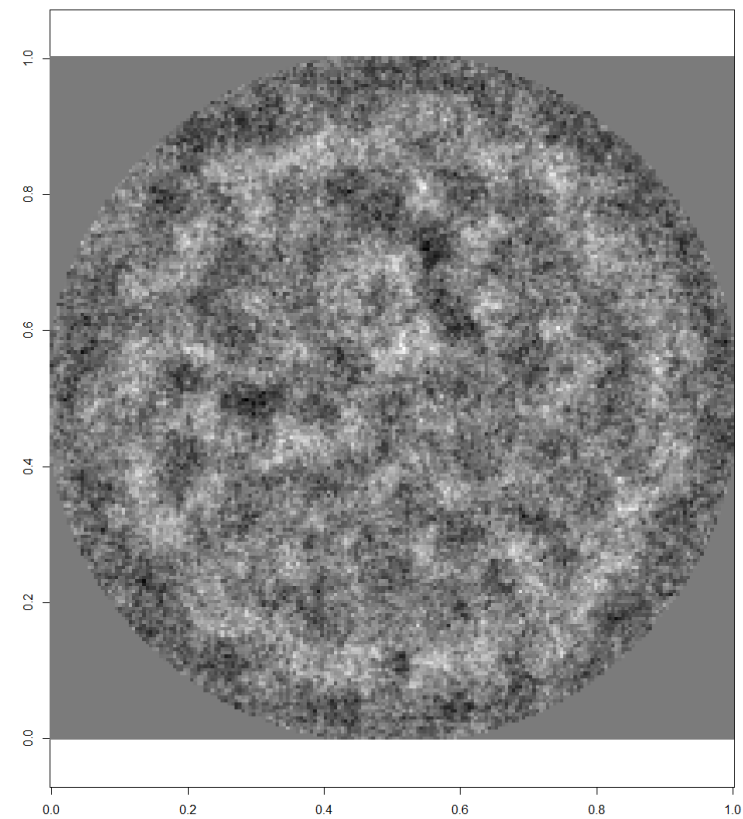
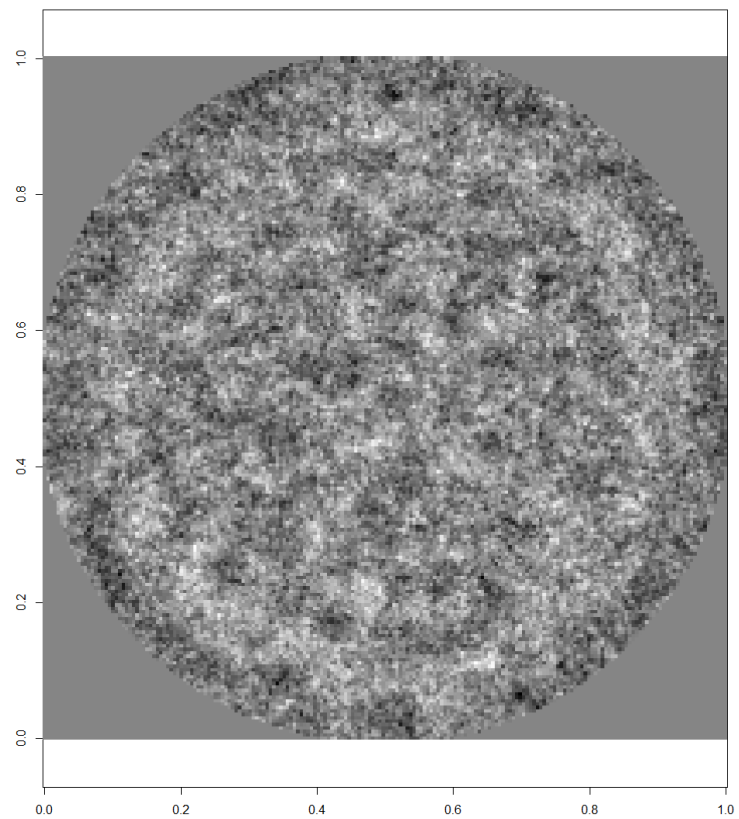


A single Cryo EM image has 40,000 pixels (shown in the previous slide). Thus each data point is a vector in 40,000 dimensional Euclidean space.

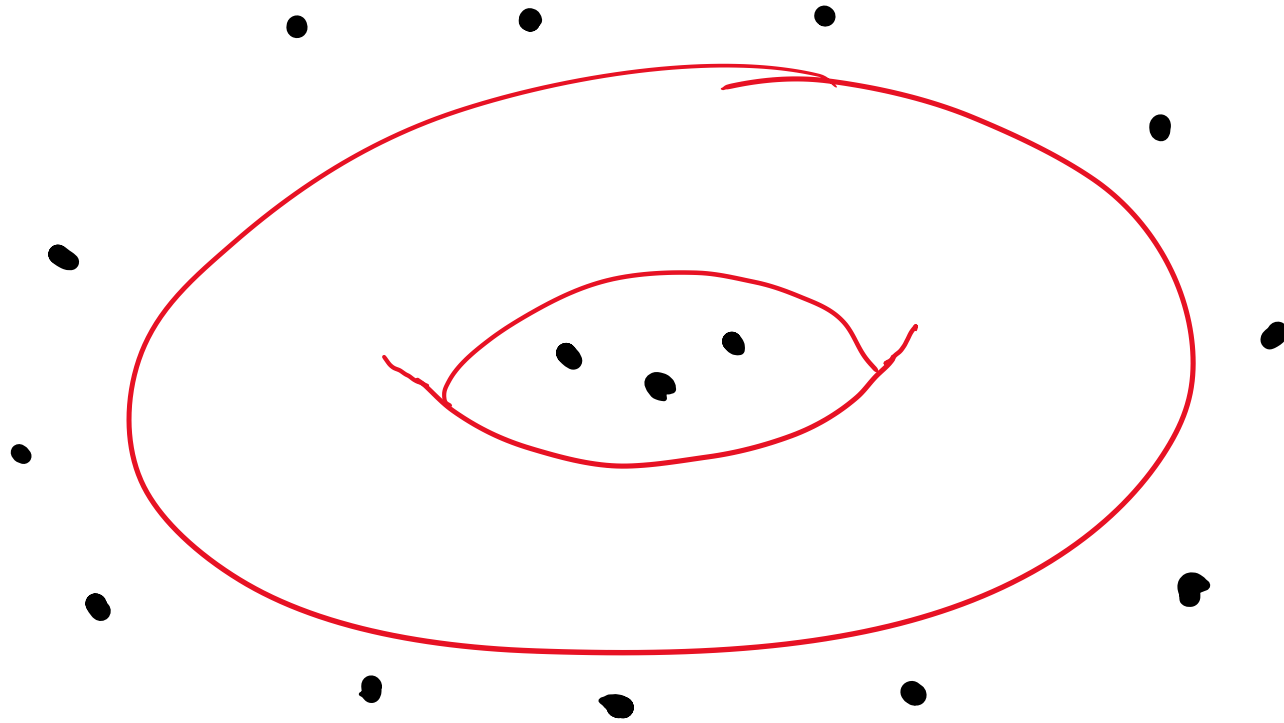
Assuming the noise is small, these points approximately lie in an orbit of SO_3 which is a 3 dimensional Lie group.

Assuming the molecule is “generic”, the orbit would be a 3-dimensional manifold diffeomorphic to SO_3 .

Typical preprocessed cryo-EM images



Data on a manifold with additive Gaussian noise



Testing the Manifold Hypothesis

- How can we test the Manifold Hypothesis?
- How many samples do we require to perform the test?
- Answers to these questions were provided in [Fefferman-Mitter-N, JAMS'16].

- We would like to infer the manifold from noisy samples.

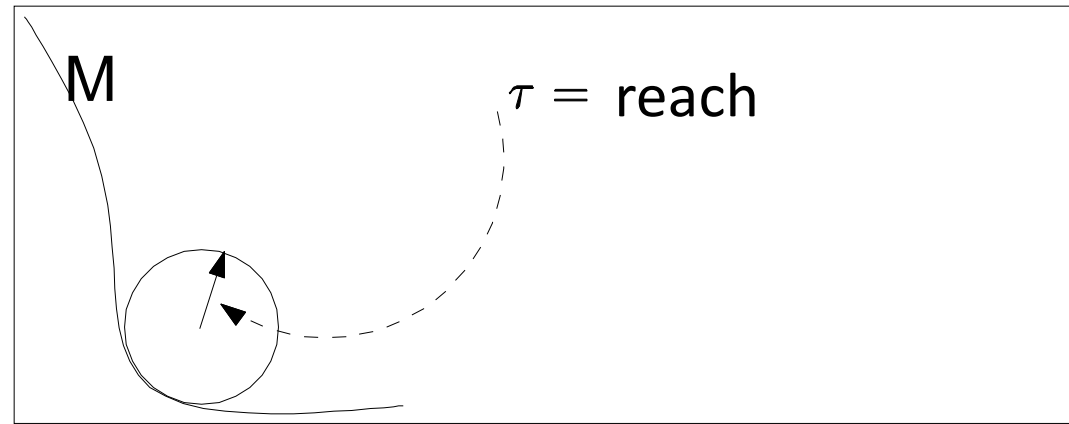
In order for this to be possible, we need to place restrictions on the manifold \mathcal{M} .

Assumptions:

1. $\mathcal{M} \subseteq \mathbb{R}^n$ has no boundary and is d -dimensional and C^2 .
2. The reach of \mathcal{M} is at least τ .
3. The d -dimensional Hausdorff measure is at most V .

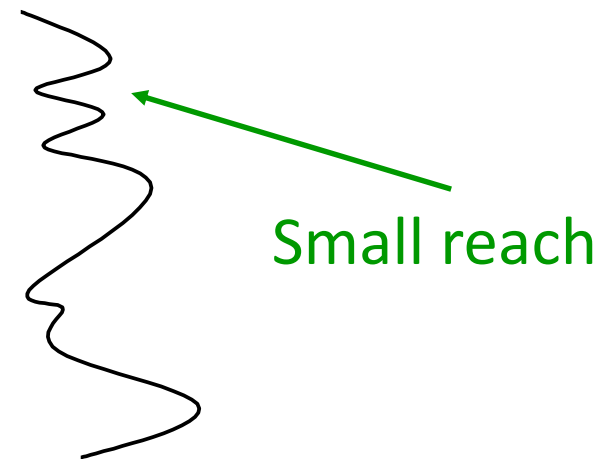
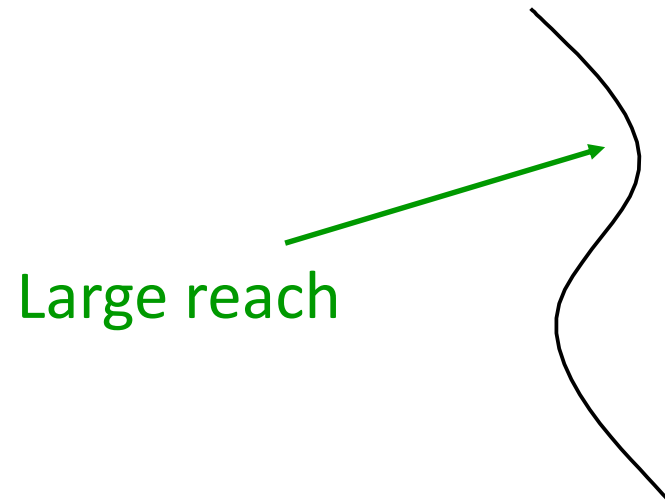
$\mathcal{M} \subseteq \mathbb{R}^n$ has
no boundary and is d –dimensional and C^2
means that for every point x in \mathcal{M}
there is $\epsilon > 0$ and a ball B_ϵ^x such that $B_\epsilon^x \cap \mathcal{M}$
is the graph of a function from a
 d –dimensional disc $Tan_x \cap B_\epsilon^x$
to the Normal space Nor_x at x .

Reach of a submanifold of \mathbb{R}^n



τ is the largest number such that for any $r < \tau$

any point at a distance r of \mathcal{M} had a unique nearest point on \mathcal{M}



For a boundaryless C^2 manifold \mathcal{M} with positive reach the d –dimensional Hausdorff measure of \mathcal{M} is equal to

$$\lim_{\epsilon \rightarrow 0} \frac{\text{vol}(\mathcal{M}_\epsilon)}{\text{vol}(B_\epsilon)},$$

where B_ϵ is the ϵ –ball of dimension $n - d$ and \mathcal{M}_ϵ is the tube of radius ϵ around \mathcal{M} .

Testing the Manifold Hypothesis [Fefferman-Mitter-N, JAMS'16]

Suppose \mathcal{P} is an unknown probability distribution supported in the **unit ball** in a separable Hilbert space, and x_1, x_2, \dots are i.i.d random samples from \mathcal{P}

Given error ϵ , dimension d , volume V , reach τ and confidence $1 - \delta$ is there an algorithm that takes a number of samples depending on these parameters and outputs whether or not there is

$$\mathcal{M} \in \mathcal{G}_\epsilon = \mathcal{G}_\epsilon(d, V, \tau)$$

such that w.p $\geq 1 - \delta$, $\int \mathbf{d}(x, \mathcal{M})^2 d\mathcal{P}(x) < \epsilon$ **?**

Sample Complexity of testing the manifold hypothesis

What is the number of samples needed for testing
the hypothesis that data lie near a low dimensional manifold?

the sample complexity of the task depends only on
the intrinsic **dimension**, **volume** and **reach**, but
not ambient dimension

Fitting manifolds

Theorem:

Let x_1, \dots, x_s be i.i.d samples from \mathcal{P} , a distribution supported on the ball of radius 1 in a separable Hilbert space. If

$$s \geq \frac{C \left(V \left(\frac{1}{\sqrt{\epsilon\tau}} + \frac{1}{\tau} \right)^{d+o(d)} + \log 1/\delta \right)}{\epsilon^2}$$

$$\text{then } \mathbb{P} \left[\sup_{\mathcal{G}_e} \left| \frac{\sum_{i=1}^s \mathbf{d}(x_i, \mathcal{M})^2}{s} - \mathbb{E}_{\mathcal{P}} \mathbf{d}(x, \mathcal{M})^2 \right| < \epsilon \right] > 1 - \delta.$$

Proof: Approximates manifolds using point clouds and uses the uniform bound for k -means.

Algorithmic question

Given N points x_1, \dots, x_N in the unit ball in \mathbb{R}^n

is there a manifold $\mathcal{M} \in \mathcal{G}_e = \mathcal{G}_e(d, CV, C^{-1}\tau)$

such that $\left(\frac{1}{N}\right) \sum_{1 \leq i \leq N} \mathbf{d}(x_i, \mathcal{M})^2 \leq C\epsilon$?

Here C is some constant depending only on d .

Theorem

There is a controlled constant C depending only on d and an Algorithm that uses linear in n but doubly exponential in d number of operations on real numbers such that given $x_1, \dots, x_N \in B_n$, with probability at least $1 - \delta$, the Algorithm outputs

1. “Yes” if there exists a manifold $\mathcal{M} \in \mathcal{G}_e(d, V, \tau)$ such that

$$\sum_{i=1}^N \mathbf{d}(x_i, \mathcal{M})^2 \leq N\epsilon,$$

2. “No” if there exists no manifold $\mathcal{M}' \in \mathcal{G}_e(d, CV, \tau/C)$ such that

$$\sum_{i=1}^N \mathbf{d}(x_i, \mathcal{M}')^2 \leq NC\epsilon,$$

Let x_1, x_2, \dots, x_N be i.i.d draws from a measure,
the logarithm of whose Radon-Nikodym derivative with respect to the
 d –dimensional Hausdorff measure on \mathcal{M} lies is L-Lipschitz.

Let ζ_1, \dots, ζ_N be a sequence of i.i.d spherical gaussians independent of x_1, \dots, x_N having a Gaussian distribution whose density at x is

$$\left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\frac{\|x\|^2}{2\sigma^2}\right),$$

where we assume

$$\frac{\tau}{Cd^C} > \sigma\sqrt{D}.$$

where D is the ambient dimension and N is the number of samples chosen and is roughly V/σ^d . We observe $y_i = x_i + \zeta_i$ for $i = 1, 2, \dots$ and wish to reconstruct \mathcal{M} up to a small error measured in Hausdorff distance.

Comparison with prior work:

Genovese et al (2014), building over Ozertem and Erdogmus (2011) provides an upper bound on the Hausdorff distance between the output manifold and the true manifold equal to $O((\frac{\log N}{N})^{\frac{2}{D+8}}) + \tilde{O}(\sigma^2 \log(\sigma^{-1}))$. Note that in order to obtain a Hausdorff distance of $c\epsilon$, one needs more than $\epsilon^{-D/2}$ samples, where D is the ambient dimension. The results of this talk guarantee (for sufficiently small σ ,) a Hausdorff distance of $Cd^7(\sigma^2 D/\tau) = O(\sigma^2)$ with less than $\frac{CV}{\omega_d(\sigma^2(D/\tau))^{d+2}} = O(\sigma^{-(2d+4)})$ samples. Genovese et al do not control the reach. We guarantee a reach of $\frac{\tau}{Cd^7}$.

Let X be a finite set of points in $E = \mathbb{R}^D$ and $X \cap B_1(x) := \{x, \tilde{x}_1, \dots, \tilde{x}_s\}$ be a set of points within a Hausdorff distance δ of some (unknown) unit d -dimensional disc $D_1(x)$ centered at x . Here $B_1(x)$ is the set of points in \mathbb{R}^D whose distance from x is less or equal to 1. We give below a simple algorithm that finds a unit d -disc centered at x within a Hausdorff distance $Cd\delta$ of $X_0 := X \cap B_1(x)$, where C is an absolute constant.

The basic idea is to choose a near orthonormal basis of d vectors from X_0 where x is taken to be the origin and let the span of this basis intersected with $B_1(x)$ be the desired disc.

Algorithm FindDisc:

1. Let x_1 be a point that minimizes $|1 - |x - x'|||$ over all $x' \in X_0$.
2. Given x_1, \dots, x_m for $m \leq d - 1$, choose x_{m+1} such that

$$\max(|1 - |x - x'|||, |\langle x_1/|x_1|, x' \rangle|, \dots, |\langle x_m/|x_m|, x' \rangle|)$$

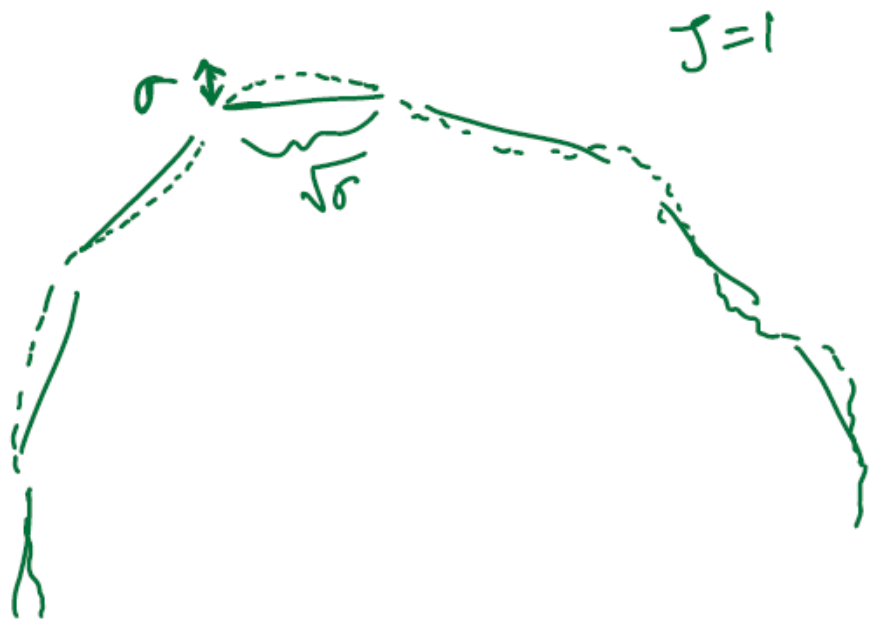
is minimized among all $x' \in X_0$ for $x' = x_{m+1}$.

Let \tilde{A}_x be the affine d -dimensional subspace containing x, x_1, \dots, x_d , and the unit d -disc $\tilde{D}_1(x)$ be $\tilde{A}_x \cap B_1(x)$.

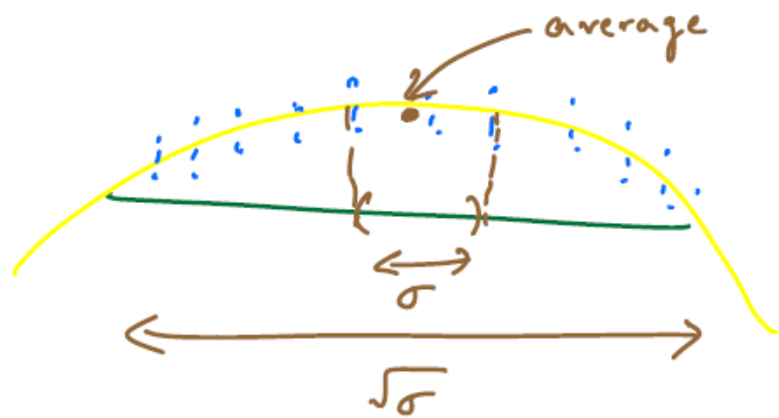
lemma: Suppose there exists a d -dimensional affine subspace A_x containing x such that $D_1(x) = A_x \cap B_1(x)$ satisfies $d_H(X_0, D_1(x)) \leq \delta$. Suppose $0 < \delta < \frac{1}{2d}$. Then $d_H(X_0, \tilde{D}_1(x)) \leq Cd\delta$, where C is an absolute constant.

We introduce a family of n dimensional balls of radius r , $\{U_i\}_{i \in [\bar{N}]}$ where the center of U_i is p_i and a family of d -dimensional embedded discs of radius r $\{D_i\}_{i \in [\bar{N}]}$, $D_i \subseteq U_i$ where D_i is centered at p_i . The D_i and the p_i are chosen by a procedure described earlier. We will need the following properties of (D_i, p_i) :

1. The Hausdorff distance between $\cup_i D_i$ and \mathcal{M} is less than $\frac{Cdr^2}{\tau} = \delta$.
2. For any $i \neq j$, $|p_i - p_j| > \frac{cr}{d}$.
3. For every $z \in \mathcal{M}$, there exists a point p_i such that $|z - p_i| < 3 \inf_{i \neq j} |p_i - p_j|$.



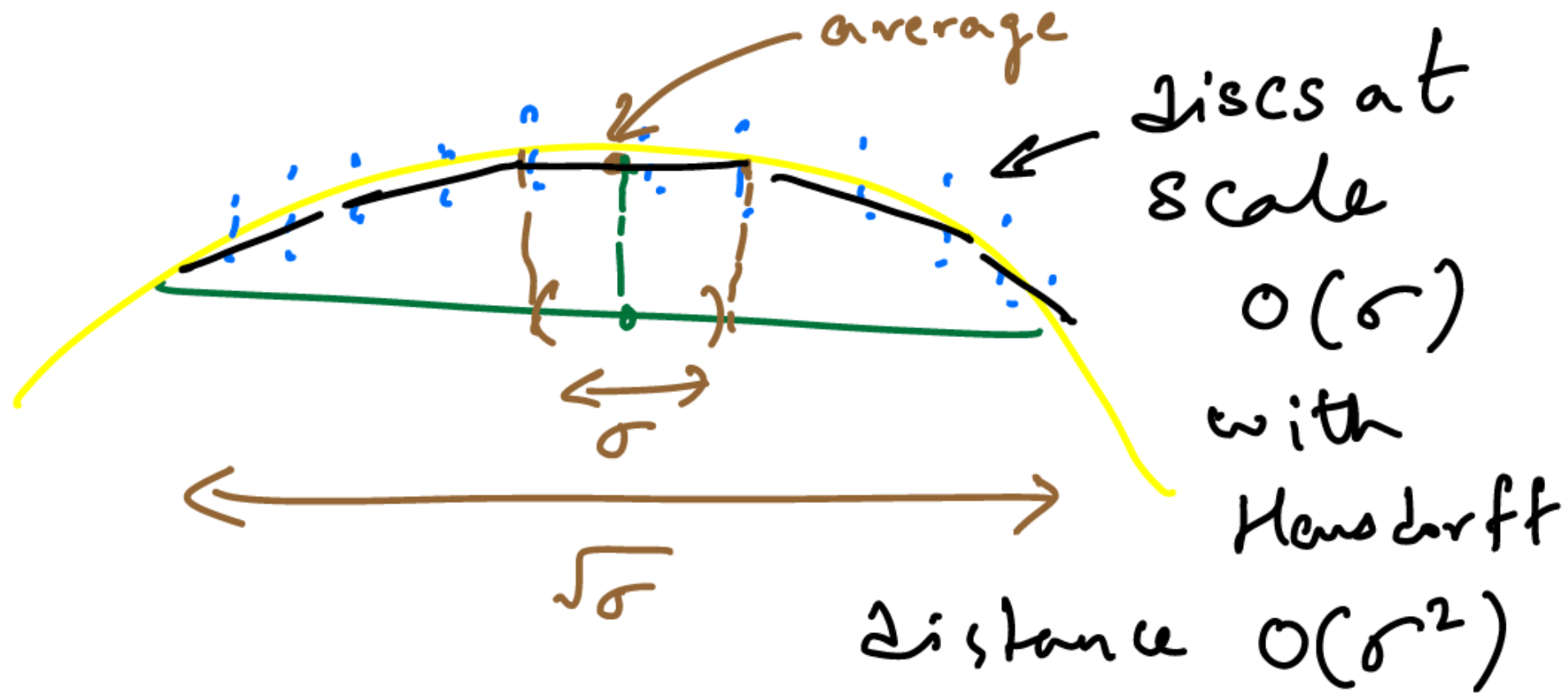
De noising by averaging:-



First, set $r = O(\sqrt{\sigma})$.
Find discs. Then using the discs as reference, at scale σ find the average of the displacements using $O(1/\sigma^{C_d})$ samples.

If the Radon-Nikodym derivative of the measure on the manifold (wrt the Hausdorff measure) is log-Lipschitz, the average is within $O(\sigma^2)$ of the average. The set of all such averages can be arranged to be within $O(\sigma^2)$ of the true manifold, in Hausdorff distance.

De noising by averaging:-



Consider the bump function $\tilde{\alpha}_i$ given by

$$\tilde{\alpha}_i(p_i + rv) = c_i(1 - \|v\|^2)^{d+2}$$

for any $v \in B_n$ and 0 otherwise. Let

$$\tilde{\alpha}(x) := \sum_i \tilde{\alpha}_i(x).$$

Let

$$\alpha_i(x) = \frac{\tilde{\alpha}_i(x)}{\sum_i \tilde{\alpha}_i(x)},$$

for each i .

Lemma:

It is possible to choose c_i such that for any z in a $\frac{r}{4d}$ neighborhood of \mathcal{M} ,

$$c^{-1} > \tilde{\alpha}(z) > c,$$

where c is a small universal constant. Further, such c_i can be computed using no more than $N_0(Cd)^{2d}$ operations involving vectors of dimension D .

Let Π^i be the orthogonal projection onto the $n - d$ -dimensional subspace containing the origin that is orthogonal to the affine span of D_i .

We define the function $F_i : U_i \rightarrow \mathbb{R}^n$ by $F_i(x) = \Pi^i(x - p_i)$. Let $\cup_i U_i = U$. We define

$$F : U \rightarrow \mathbb{R}^n$$

by $F(x) = \sum_i \alpha_i(x) F_i(x)$.

Given a symmetric matrix A such that A has $n - d$ eigenvalues in $(1/2, 3/2)$ and d eigenvalues in $(-1/2, 1/2)$, let $\Pi_{hi}(A)$ denote the projection onto the span of the eigenvectors corresponding to the top $n - d$ eigenvalues.

For $x \in \cup_i U_i$, we define $\Pi_x = \Pi_{hi}(A_x)$ where $A_x = \sum_i \alpha_i(x) \Pi^i$. Let \tilde{U}_i be defined as the $\frac{cr}{d}$ -Euclidean neighborhood of D_i inside U_i . Note that Π_x is C^2 when restricted to $\cup_i \tilde{U}_i$, because the $\alpha_i(x)$ are C^2 and when x is in this set, $c < \sum_i \tilde{\alpha}_i(x) < c^{-1}$, and for any i, j such that $\alpha_i(x) \neq 0 \neq \alpha_j(x)$, we have $\|\Pi^i - \Pi^j\|_F < Cd\delta$.

We define the output manifold \mathcal{M}_o to be the set of all points x such that $x \in \tilde{U}_i$ for some i and $\Pi_x F(x) = 0$.

Recall the bump function $\tilde{\alpha}_i$ given by

$$\tilde{\alpha}_i(p_i + rv) = c_i(1 - \|v\|^2)^{d+2}$$

for any $v \in B_n$ and 0 otherwise.

Observe that

$$\sum_i |\partial_v \tilde{\alpha}_i(x)|^{\frac{d+2}{d+1}} \leq Cd \|(\tilde{\alpha}_i(x))_i\|_1 \leq Cd.$$

Recall

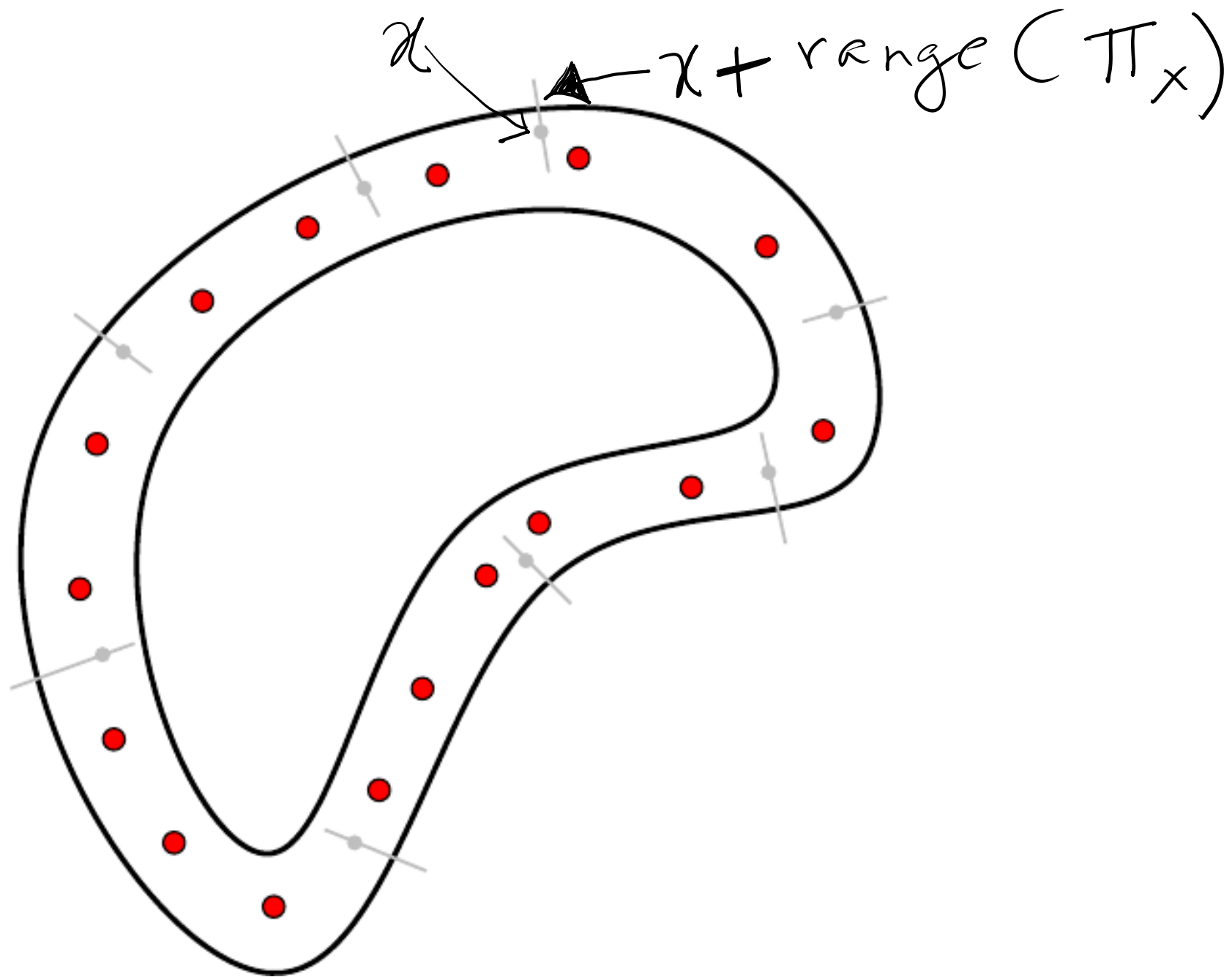
$$\tilde{\alpha}(x) := \sum_i \tilde{\alpha}_i(x).$$

Thus by the Hölder inequality,

$$|\partial_v \tilde{\alpha}| \leq \|\partial_v \tilde{\alpha}_i(x)\|_{\frac{d+2}{d+1}} (\mathbf{1})_{d+2} \leq Cd^2.$$

Recall that for $x \in \cup_i U_i$, we define $\Pi_x = \Pi_{hi}(A_x)$ where $A_x = \sum_i \alpha_i(x) \Pi^i$. Let \tilde{U}_i be defined as the $\frac{cr}{d}$ -Euclidean neighborhood of D_i inside U_i . Note that Π_x is C^2 when restricted to $\cup_i \tilde{U}_i$, because the $\alpha_i(x)$ are C^2 and when x is in this set, $c < \sum_i \tilde{\alpha}_i(x) < c^{-1}$, and for any i, j such that $\alpha_i(x) \neq 0 \neq \alpha_j(x)$, we have $\|\Pi^i - \Pi^j\|_F < Cd\delta$.

We define the output manifold \mathcal{M}_o to be the set of all points x such that $x \in \tilde{U}_i$ for some i and $\Pi_x F(x) = 0$.



We define the output manifold \mathcal{M}_o to be the set of all points x such that $x \in \tilde{U}_i$ for some i and $\Pi_x F(x) = 0$.

Restricted to \tilde{U}_i , if $\Pi_x F(x) \neq 0$, then $\Pi^i \Pi_x F(x) \neq 0$.

Let Π^i be the orthogonal projection onto the $n - d$ -dimensional subspace containing the origin that is orthogonal to the affine span of D_i .

We define the function $F_i : U_i \rightarrow \mathbb{R}^n$ by $F_i(x) = \Pi^i(x - p_i)$. Let $\cup_i U_i = U$. We define

$$F : U \rightarrow \mathbb{R}^n$$

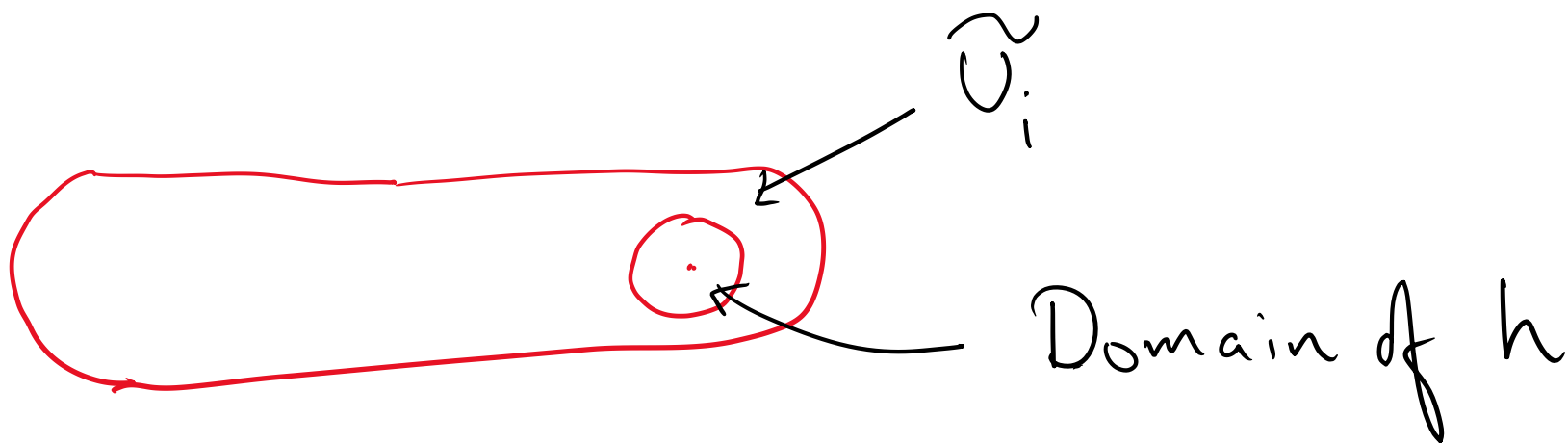
by $F(x) = \sum_i \alpha_i(x) F_i(x)$.

Given a symmetric matrix A such that A has $n - d$ eigenvalues in $(1/2, 3/2)$ and d eigenvalues in $(-1/2, 1/2)$, let $\Pi_{hi}(A)$ denote the projection onto the span of the eigenvectors corresponding to the top $n - d$ eigenvalues.

Let T denote a translation composed with a scaling, mapping from a ball of radius 1 around the origin to a ball of radius $\frac{cr}{d}$ around x_0 , contained in \tilde{U}_i .

Let

$$h = \Pi_i \Pi_x F \circ T.$$



Quantitative implicit function Theorem:

Let $h : \mathbb{R}^{m+n} \rightarrow \mathbb{R}^n$ be a C^2 -function, $h : (x, y) \mapsto h(x, y)$. Let $g : B_{m+n} \rightarrow \mathbb{R}^{m+n}$ be defined by $g : (x, y) \mapsto (x, h(x, y))$.

Suppose the Jacobian of g , Jac_g satisfies $\|Jac_g - I\| < \epsilon^2/4$ on B_{m+n} and that for any vector $v \in \mathbb{R}^{m+n}$,

$$\left\| \frac{\partial^2 g(x)}{\partial v^2} \right\| \leq \left(\frac{\epsilon^2}{4} \right) \|v\|^2$$

where $\epsilon \in [0, 1]$. Suppose also that $\|g(0)\| < \frac{\epsilon^2}{20}$ for the same ϵ .

Quantitative implicit function Theorem:

On the domain of definition of f , $g(B_{m+n})$

$$f((x, y)) = (x, e(x, y))$$

for an appropriate e and in particular, for $\|x\| \leq \frac{\eta}{2}$, where $\eta \in [0, 1]$,

$$f((x, 0)) = (x, e(x, 0))$$

and

$$\|(x, e(x, 0))\| \leq \frac{8}{5} \left(\frac{\epsilon^2}{20} + \frac{\eta}{2} \right).$$

Finally, for any $w \in \mathbb{R}^n$ such that $\|w\| = 1$, $\|Hess(e \cdot w)\| \leq \frac{16\epsilon^2}{(4-\epsilon)^3}$.

Theorem: Suppose $C\sigma\sqrt{D}$ is less than $\frac{\tau}{Cd^C}$. The reach of \mathcal{M}_o is at least $\frac{1}{Cd^7}\tau$ and the Hausdorff distance between \mathcal{M}_o and \mathcal{M} is less or equal to

$$Cd^7(D\frac{\sigma^2}{\tau}).$$

Proof: Use Cauchy's Integral formula, to write

$$\Pi_x = \frac{1}{2\pi\iota} \oint_{\gamma} (zI - A_x)^{-1} dz,$$

for suitable γ . Use Hölder Inequalities to get good bounds on the first and second derivatives of $\Pi_x F(x)$ and then apply a dimension-free quantitative form of the implicit function theorem.

Conclusion

- We bring down the Hausdorff distance between the true and output manifolds to $O(\sigma^2 D)$ while bringing down the sample and computational complexities to depend only on the intrinsic dimension.

Thank You!