

# Topological consistency via kernel estimation

Amoukou Salim

Statistics

May 31, 2019

The goal is to estimate the topology of a set  $D_L = \{x \in \mathbb{R}^d : f(x) \geq L\}$

- Why: Clustering, pattern recognition, econometric...
- With:  $D_n = \{(X_1, Y_1), \dots (X_n, Y_n)\}$
- Simple idea: Use the plug-in estimators.

The goal is to estimate the topology of a set  $D_L = \{x \in \mathbb{R}^d : f(x) \geq L\}$

- Why: Clustering, pattern recognition, econometric...
- With:  $D_n = \{(X_1, Y_1), \dots (X_n, Y_n)\}$
- Simple idea: Use the plug-in estimators.

The goal is to estimate the topology of a set  $D_L = \{x \in \mathbb{R}^d : f(x) \geq L\}$

- Why: Clustering, pattern recognition, econometric...
- With:  $D_n = \{(X_1, Y_1), \dots (X_n, Y_n)\}$
- Simple idea: Use the plug-in estimators.

## Theorem

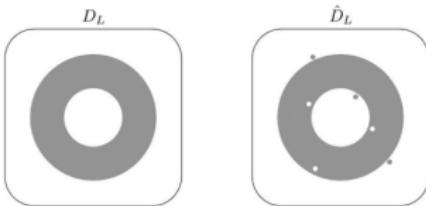
Let  $f$  be a density with compact support  $S$ ,  $f_n$  a sequence of density estimators and  $\hat{S} = \{f_n > \alpha_n\}$  where  $\alpha_n \searrow 0$ .

If we assume that:

- $\lambda(E_0) = 0$  where  $E_0 = \{x \in S | f(x) = 0\}$
- $\alpha_n^{-1} \int |f_n - f| d\lambda \rightarrow 0$  (as)

then  $d_\lambda(S, \hat{S}) \rightarrow 0$

**Goal:** Good reconstruction according to Hausdorff distance  $\not\Rightarrow$  Good recovery of the topology of  $X$ .



# How we describe the topology of $D_L$ ?



Sorbonne University

The group of homology of  $X$  : A set of vector space  $\{H_k(X)\}_{k=0}^{\infty}$   
→ **Intuitively**:

- $\text{rank}(H_0(X)) \cong$  connected components
- $\text{rank}(H_k(X)) \cong k\text{-cycle of } X$

The persistent homology of  $X$ :  $PH_*(X)$

- track the evolution of the homology as  $L$  decreasing or increasing.
- bottleneck distance to compare different persistent diagram:  
$$d_B(PH_*(X), PH_*(Y)) = \inf_{\gamma} \sup_{p \in Dgm(X)} \|p - \gamma(p)\|$$

# Naive approach



Sorbonne University

Naive approach to recover the homology of an unknown space  $S$  from a random sample  $\mathcal{X}$ :  $U(\mathcal{X}, r) = \bigcup_{x \in \mathcal{X}} B_r(x)$  for some choice of  $r$ .

## Theorem

Let  $\mu$  a probability distribution and  $K = \text{supp}(\mu)$  st:

- $\text{reach}(K) > 0$
- $\mu$  is  $(a, b)$  standard at scale  $r_0 > 0$

Let  $X_1, \dots, X_n$  be iid sample of  $\mu$  and  $\epsilon < \frac{1}{17} \text{reach}(K) \wedge 2r_0$  then  
 $\forall n \geq C_{a,b} \log(\frac{1}{\epsilon}) + \log(\frac{1}{\epsilon})/\epsilon^b$

$$\mathbb{P}[K^{\frac{1}{2}\text{reach}(K)} \text{ and } \mathbb{X}_n^{4\epsilon} \text{ are homotopy equivalent}] \geq 1 - \eta$$

Pb: Strong assumption on the density and the shape of the level set.

Our solution: We will use the same technique but with weak assumption to recover the homology.

# Naive approach



Sorbonne University

Naive approach to recover the homology of an unknown space  $S$  from a random sample  $\mathcal{X}$ :  $U(\mathcal{X}, r) = \bigcup_{x \in \mathcal{X}} B_r(x)$  for some choice of  $r$ .

## Theorem

Let  $\mu$  a probability distribution and  $K = \text{supp}(\mu)$  st:

- $\text{reach}(K) > 0$
- $\mu$  is  $(a, b)$  standard at scale  $r_0 > 0$

Let  $X_1, \dots, X_n$  be iid sample of  $\mu$  and  $\epsilon < \frac{1}{17} \text{reach}(K) \wedge 2r_0$  then  
 $\forall n \geq C_{a,b} \log(\frac{1}{\epsilon}) + \log(\frac{1}{\epsilon})/\epsilon^b$

$\mathbb{P}[K^{\frac{1}{2}\text{reach}(K)} \text{ and } \mathbb{X}_n^{4\epsilon} \text{ are homotopy equivalent}] \geq 1 - \eta$

Pb: Strong assumption on the density and the shape of the level set.

Our solution: We will use the same technique but with weak assumption to recover the homology.

# Naive approach



Sorbonne University

Naive approach to recover the homology of an unknown space  $S$  from a random sample  $\mathcal{X}$ :  $U(\mathcal{X}, r) = \bigcup_{x \in \mathcal{X}} B_r(x)$  for some choice of  $r$ .

## Theorem

Let  $\mu$  a probability distribution and  $K = \text{supp}(\mu)$  st:

- $\text{reach}(K) > 0$
- $\mu$  is  $(a, b)$  standard at scale  $r_0 > 0$

Let  $X_1, \dots, X_n$  be iid sample of  $\mu$  and  $\epsilon < \frac{1}{17} \text{reach}(K) \wedge 2r_0$  then  
 $\forall n \geq C_{a,b} \log(\frac{1}{\epsilon}) + \log(\frac{1}{\epsilon})/\epsilon^b$

$$\mathbb{P}[K^{\frac{1}{2}\text{reach}(K)} \text{ and } \mathbb{X}_n^{4\epsilon} \text{ are homotopy equivalent}] \geq 1 - \eta$$

Pb: Strong assumption on the density and the shape of the level set.

Our solution: We will use the same technique but with weak assumption to recover the homology.

# Naive approach



Sorbonne University

Naive approach to recover the homology of an unknown space  $S$  from a random sample  $\mathcal{X}$ :  $U(\mathcal{X}, r) = \bigcup_{x \in \mathcal{X}} B_r(x)$  for some choice of  $r$ .

## Theorem

Let  $\mu$  a probability distribution and  $K = \text{supp}(\mu)$  st:

- $\text{reach}(K) > 0$
- $\mu$  is  $(a, b)$  standard at scale  $r_0 > 0$

Let  $X_1, \dots, X_n$  be iid sample of  $\mu$  and  $\epsilon < \frac{1}{17} \text{reach}(K) \wedge 2r_0$  then  
 $\forall n \geq C_{a,b} \log(\frac{1}{\epsilon}) + \log(\frac{1}{\epsilon})/\epsilon^b$

$$\mathbb{P}[K^{\frac{1}{2}\text{reach}(K)} \text{ and } \mathbb{X}_n^{4\epsilon} \text{ are homotopy equivalent}] \geq 1 - \eta$$

Pb: Strong assumption on the density and the shape of the level set.

Our solution: We will use the same technique but with weak assumption to recover the homology.

The main assumption on  $f$  is "tameless". Let  $D_L$  define like before.

## Definition

- $L$  is called homological regular value if  $\exists \epsilon > 0$  st  $\forall v_1 \leq v_2$  in  $[L - 2\epsilon, L + 2\epsilon]$ , the inclusion map  $H_*(D_{v1}) \hookrightarrow H_*(D_{v2})$  is an isomorphism. Otherwise, we say that  $L$  is **homological critical value**
- A function is called tame if it has a **finite** homological critical value and  $\text{rank}(H_*(D_L))$  is finite  $\forall L$

# To sum up



Our goal is to present a consistent method for recovering the homology of a given set  $D_L$ . We will examine two cases :

- Density function  $p$  with kernel estimator  $f_n(x) = \frac{1}{C_k nr^d} \sum K_r(x - X_i)$ 
  - $p_{max} = \sup_{x \in \mathbb{R}^d} p(x) < +\infty$
  - $D_L$  is bounded for all  $L > 0$
- Regression function  $f(x) = \mathbb{E}(Y|X = x)$  with kernel estimator  $f_n(x) = \frac{\sum Y_i K_r(x - X_i)}{\sum K_r(x - X_i)}$ , if  $p$  is the marginal density of  $X$  then
  - $\text{supp}(p)$  compact and  $p_{min} = \inf_{x \in \mathbb{R}^d} p(x) > 0$
  - $|Y_i| < Y_{max}$  almost surely
- Kernel estimator satisfies the following:
  - $\text{supp}(K) \subset B_1(0)$
  - $K(x) \in [0, 1]$ , and  $K(0) = 1$
  - $\int K(x)dx = C_K$  for  $C_K \in (0, 1)$

# To sum up



Sorbonne University

Our goal is to present a consistent method for recovering the homology of a given set  $D_L$ . We will examine two cases :

- Density function  $p$  with kernel estimator  $f_n(x) = \frac{1}{C_k nr^d} \sum K_r(x - X_i)$ 
  - $p_{max} = \sup_{x \in \mathbb{R}^d} p(x) < +\infty$
  - $D_L$  is bounded for all  $L > 0$
- Regression function  $f(x) = \mathbb{E}(Y|X = x)$  with kernel estimator  $f_n(x) = \frac{\sum Y_i K_r(x - X_i)}{\sum K_r(x - X_i)}$ , if  $p$  is the marginal density of  $X$  then
  - $\text{supp}(p)$  compact and  $p_{min} = \inf_{x \in \mathbb{R}^d} p(x) > 0$
  - $|Y_i| < Y_{max}$  almost surely
- Kernel estimator satisfies the following:
  - $\text{supp}(K) \subset B_1(0)$
  - $K(x) \in [0, 1]$ , and  $K(0) = 1$
  - $\int K(x)dx = C_K$  for  $C_K \in (0, 1)$

# To sum up



Our goal is to present a consistent method for recovering the homology of a given set  $D_L$ . We will examine two cases :

- Density function  $p$  with kernel estimator  $f_n(x) = \frac{1}{C_k nr^d} \sum K_r(x - X_i)$ 
  - $p_{max} = \sup_{x \in \mathbb{R}^d} p(x) < +\infty$
  - $D_L$  is bounded for all  $L > 0$
- Regression function  $f(x) = \mathbb{E}(Y|X = x)$  with kernel estimator  $f_n(x) = \frac{\sum Y_i K_r(x - X_i)}{\sum K_r(x - X_i)}$ , if  $p$  is the marginal density of  $X$  then
  - $\text{supp}(p)$  compact and  $p_{min} = \inf_{x \in \mathbb{R}^d} p(x) > 0$
  - $|Y_i| < Y_{max}$  almost surely
- Kernel estimator satisfies the following:
  - $\text{supp}(K) \subset B_1(0)$
  - $K(x) \in [0, 1]$ , and  $K(0) = 1$
  - $\int K(x)dx = C_K$  for  $C_K \in (0, 1)$

# To sum up



Our goal is to present a consistent method for recovering the homology of a given set  $D_L$ . We will examine two cases :

- Density function  $p$  with kernel estimator  $f_n(x) = \frac{1}{C_k nr^d} \sum K_r(x - X_i)$ 
  - $p_{max} = \sup_{x \in \mathbb{R}^d} p(x) < +\infty$
  - $D_L$  is bounded for all  $L > 0$
- Regression function  $f(x) = \mathbb{E}(Y|X = x)$  with kernel estimator  $f_n(x) = \frac{\sum Y_i \sum K_r(x - X_i)}{\sum K_r(x - X_i)}$ , if  $p$  is the marginal density of  $X$  then
  - $\text{supp}(p)$  compact and  $p_{min} = \inf_{x \in \mathbb{R}^d} p(x) > 0$
  - $|Y_i| < Y_{max}$  almost surely
- Kernel estimator satisfies the following:
  - $\text{supp}(K) \subset B_1(0)$
  - $K(x) \in [0, 1]$ , and  $K(0) = 1$
  - $\int K(x)dx = C_K$  for  $C_K \in (0, 1)$

# First procedure: Naive estimator



Sorbonne University

- 1 Use the dataset to construct an Kernel estimator  $\hat{f}$
- 2 Using the estimator  $\hat{f}$ , define  $\mathcal{X}^L = \{X_i : \hat{f}(X_i) \geq L\}$
- 3 Consider  $U(\mathcal{X}, r)$  as an estimate of  $D_L$ , and  $\hat{D}_L(n, r) = H_*(U(\mathcal{X}, r))$  as an estimate of  $H_*(D_L)$

Pb :  $\hat{f}$  may occur error in the step 2

Solution : Use filtering step like in the course to overcome the nosiness of  $\hat{D}_L(n, r)$ .

# First procedure: Naive estimator



Sorbonne University

- 1 Use the dataset to construct an Kernel estimator  $\hat{f}$
- 2 Using the estimator  $\hat{f}$ , define  $\mathcal{X}^L = \{X_i : \hat{f}(X_i) \geq L\}$
- 3 Consider  $U(\mathcal{X}, r)$  as an estimate of  $D_L$ , and  $\hat{D}_L(n, r) = H_*(U(\mathcal{X}, r))$  as an estimate of  $H_*(D_L)$

Pb :  $\hat{f}$  may occur error in the step 2

Solution : Use filtering step like in the course to overcome the nosiness of  $\hat{D}_L(n, r)$ .

# First procedure: Naive estimator



Sorbonne University

- 1 Use the dataset to construct an Kernel estimator  $\hat{f}$
- 2 Using the estimator  $\hat{f}$ , define  $\mathcal{X}^L = \{X_i : \hat{f}(X_i) \geq L\}$
- 3 Consider  $U(\mathcal{X}, r)$  as an estimate of  $D_L$ , and  $\hat{D}_L(n, r) = H_*(U(\mathcal{X}, r))$  as an estimate of  $H_*(D_L)$

Pb :  $\hat{f}$  may occur error in the step 2

**Solution :** Use filtering step like in the course to overcome the nosiness of  $\hat{D}_L(n, r)$ .

# Use image to filter



Since we have  $\hat{D}_{L+\epsilon}(n, r) \subset \hat{D}_{L-\epsilon}(n, r)$ , we have the homology map :

$$i_* : H_*(\hat{D}_{L+\epsilon}(n, r)) \hookrightarrow H_*(\hat{D}_{L+\epsilon}(n, r))$$

The new estimator :  $\hat{H}(L, \epsilon, n) = \text{Im}(i_*)$

$$\begin{array}{ccccc} D_{L+2\epsilon} & \searrow & D_L & \searrow & D_{L-2\epsilon} \\ & \hat{D}_{L+\epsilon}(n, r) & \xrightarrow{i_*} & \hat{D}_{L-\epsilon}(n, r) & \end{array},$$

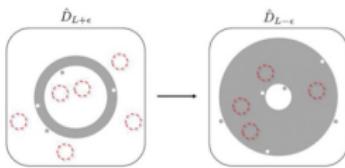


Figure: The intuition behind

To have the precedent inclusion sequence, we require the following regularity condition on  $L$ .

## Definition

$L$  is  $\epsilon$ -regular if :

$$\partial D_{L+2\epsilon} \cap \partial D_{L+\frac{3}{2}\epsilon} = D_{L+\frac{1}{2}\epsilon} \cap \partial D_L = \partial D_L \cap D_{L-\frac{1}{2}\epsilon} = \partial D_{L-\frac{3}{2}\epsilon} \cap \partial D_{L-2\epsilon} = \emptyset$$

*In particular, if  $f$  is continuous in  $f^{-1}([L - 2\epsilon, L + 2\epsilon])$  so  $L$  is  $\epsilon$ -regular.*

To have the precedent inclusion sequence, we require the following regularity condition on  $L$ .

## Definition

$L$  is  $\epsilon$ -regular if :

$$\partial D_{L+2\epsilon} \cap \partial D_{L+\frac{3}{2}\epsilon} = D_{L+\frac{1}{2}\epsilon} \cap \partial D_L = \partial D_L \cap D_{L-\frac{1}{2}\epsilon} = \partial D_{L-\frac{3}{2}\epsilon} \cap \partial D_{L-2\epsilon} = \emptyset$$

*In particular, if  $f$  is continuous in  $f^{-1}([L - 2\epsilon, L + 2\epsilon])$  so  $L$  is  $\epsilon$ -regular.*

# The main theorem of the paper



Sorbonne University

## Theorem

Let  $L > 0$  and  $\epsilon \in (0, \frac{L}{2})$  be such that  $f(x)$  has no critical value in the range  $[L - 2\epsilon, L + 2\epsilon]$ , if  $r \rightarrow 0$  and  $nr^d \rightarrow +\infty$  then:

$$\mathbb{P}(H_*(L, \epsilon, n)) \cong H_*(D_L)) \geq 1 - 6ne^{-C_{\epsilon/2}^* nr^d}$$

In particular, if  $nr^d \geq D \log(n)$  with  $D > (C_{\epsilon/2})^{-1}$

$$\mathbb{P}(H_*(L, \epsilon, n)) \cong H_*(D_L)) = 1$$

To prove the theorem, we need 3 lemmas. The first is :

## Lemma

For every  $L > 0$ , and  $\epsilon \in (0, L)$ , if  $r \rightarrow 0$  and  $nr^d \rightarrow \infty$ , then there exists a constant  $C_\epsilon$  such that for  $n$  large enough we have :

$$\mathbb{P}(\exists X_i \notin D_{L-\epsilon}^\uparrow(r) : f_n(X_i) \geq L) \leq ne^{-C_\epsilon^* nr^d} \quad (1)$$

$$\mathbb{P}(\exists X_i \in D_{L-\epsilon}^\uparrow(r) : f_n(X_i) \leq L) \leq ne^{-C_\epsilon^* nr^d} \quad (2)$$

# Proof lemma 1



Sorbonne University

Recall that:  $\text{Supp}(K) \subset B_1(0)$ ,  $K(r) \in [0, 1]$  and  $K(0)=1$

The kernel estimator is:  $\hat{f}_m(\alpha) = \frac{\sum K_p(\alpha - x_i)}{C_K m^{1/d}}$

We have:

$$\begin{aligned} \mathbb{P}\left(\exists x_i \notin B_{L-\varepsilon}^c(n) : \hat{f}_m(x_i) \geq L\right) &\leq m \mathbb{P}\left(X_i \in \left(B_{L-\varepsilon}^c(n)\right)^c : \hat{f}_m(x_i) \geq L\right) \\ &\leq m \int_{\left(B_{L-\varepsilon}^c(n)\right)^c} \hat{f}_m(\alpha) \boxed{\mathbb{P}(\hat{f}_m(x_i) \geq L | x=\alpha)} \end{aligned}$$

$$\begin{aligned} \boxed{\mathbb{P}(\hat{f}_m(x_i) \geq L | x=\alpha)} &= \mathbb{P}\left(K_p(0) + \sum_{i=2}^m K_p(\alpha - x_i) \geq L C_K m^{1/d}\right) \\ &= \mathbb{P}\left(\sum_{i=2}^m Z_i \geq m(L C_K m^{1/d} - p_p(\alpha)) + p_p(\alpha) - 1\right) \end{aligned}$$

Where  $Z_i = K_p(\alpha - x_i) - p_p(\alpha)$  and  $\boxed{p_p(\alpha) = \mathbb{E} \sum K_p(\alpha - x_i)}$

$$\boxed{p_p(\alpha) = \mathbb{E} \sum K_p(\alpha - x_i)} = \int_{B_p(\alpha)} f_X(\xi) K_p(\xi - x_i) d\xi \underset{\substack{\uparrow \\ \alpha \notin B_{L-\varepsilon}^c(n)}}{\leq} (L - \varepsilon) C_K m^{1/d}$$

# Proof lemma 1



So reassemble the puzzle and we have:

$$\boxed{\Pr(\hat{f}_{\hat{b}_m}(x_i) \geq L | x = \alpha)} \leq \Pr\left(\sum_{i=2}^m z_i \geq \varepsilon C_{K, m, d} - 1\right)$$

Note that  $|z_i| \leq 1$  and  $\text{Var}(z_i) = \mathbb{E}[z_i^2] - (\mathbb{E}[z_i])^2 \leq p_{\max} C_{K, m, d}$

Using Bernstein inequality with  $t = \varepsilon C_{K, m, d} - L$  and the limit we have:

$$\boxed{\Pr(\hat{f}_{\hat{b}_m}(x_i) \geq L | x = \alpha)} \leq e^{-\frac{\varepsilon^2}{C} m d}$$

$$\text{where } C^* = \frac{\varepsilon^2 C_K}{3p_{\max} + \varepsilon}$$

- We do same technique for the second and 3rd Regression case.

## Lemma

For every  $L > 0$ , and  $\epsilon \in (0, L)$ , if  $r \rightarrow 0$  and  $nr^d \rightarrow \infty$ , then there exists a constant  $C_\epsilon$  such that for  $n$  large enough we have :

$$\mathbb{P}(D_{L+\epsilon}^\downarrow(2r) \subset \hat{D}_L(n, r) \subset D_{L-\epsilon}^\uparrow(2r)) \geq 1 - 3ne^{-C_\epsilon nr^d} \quad (3)$$

# Proof lemma 2



- Recall that we want  $\text{IP}(\mathbb{B}_{L-\epsilon}^b(2n) \subset \mathbb{B}_L^b(m, n) \subset \mathbb{B}_{L-\epsilon}^b(2n)) \geq 1 - 3m^{-\frac{C_2 m n d}{C_1 r}}$

$$- \text{IP}(\mathbb{B}_L^b(m, n) \not\subset \mathbb{B}_{L-\epsilon}^b(2n)) \leq \text{IP}(\exists x_i \notin \mathbb{B}_{L-\epsilon}^b(n) : f_m(x_i) \geq L) \leq m^{-\frac{C_2 m n d}{C_1 r}}$$

For the second inclusion, we have to make that:

- Since  $\mathbb{B}_{L-\epsilon}^b(2n)$  is bounded, let  $S \subset \mathbb{B}_{L-\epsilon}^b(2n)$  be a  $\delta n$ -covering.
- If  $\exists x \in \mathbb{B}_{L-\epsilon}^b(2n)$  that is not covered by  $\cup(S_n)_n$  so  $\exists s \in S$  st  $s$  is not covered by  $\cup(S_m)(1-\delta)n$ .

Lemma:  $\forall y \in \mathbb{R}^L : \|x - y\| > n$ , and let  $s \in S$  such that  $\|x - s\| \leq \delta n \Rightarrow$

$$\|y - s\| = \|y - x + x - s\| \geq \|y - x\| - \|x - s\|$$

$$\geq n - \delta n \geq n(1 - \delta)$$

Hence on pent dire:

$$\begin{aligned} & \text{IP}(\mathbb{B}_{L-\epsilon}^b(2n) \subset \mathbb{B}_L^b(m, n)) \leq \text{IP}(\exists s \in S : B_{(-\delta)n}^b(s) \cap S_m^L = \emptyset) \\ & \leq \text{IP}(\exists s \in S : B_{(-\delta)n}^b(s) \cap S_m^L = \emptyset ; \mathbb{B}_{L-\epsilon}^b(n) \cap S_m^L \neq \emptyset) + \text{IP}(\exists s \in S : B_{(-\delta)n}^b(s) \cap S_m^L = \emptyset ; \mathbb{B}_{L-\epsilon}^b(n) \cap S_m^L = \emptyset) \\ & \leq \text{IP}(\exists s \in S : B_{(-\delta)n}^b(s) \cap S_m^L = \emptyset) + \text{IP}(\mathbb{B}_{L-\epsilon}^b(n) \cap S_m^L \neq \emptyset) \end{aligned}$$

Original sample does not

$$\leq C_1 r^{-d} e^{-C_2 m n d}$$

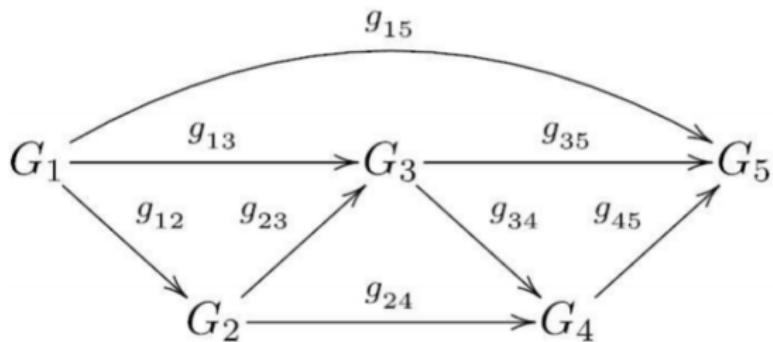
$O(n)$

get red to many point

$$\leq m^{-\frac{C_2 m n d}{C_1 r}} \Rightarrow \leq 2 m^{-\frac{C_2 m n d}{C_1 r}}$$

## Lemma 3

Consider the following commutative diagram of group :



Define  $G_{i,j} = \text{Im}(g_{i,j}) \subset G_j$ .

If  $g_{3,5} : G_3 \rightarrow G_5$  is an isomorphism from  $G_3$  to  $G_{1,5}$ . Then the map  $g_{3,4} : G_3 \rightarrow G_4$  is an isomorphism from  $G_3$  to  $G_{2,4}$ . In particular, we have

$$G_3 \cong G_{2,4}$$

# The proof of the theorem



Sorbonne University

## Proof.

Recall that we want to prove that: if  $r \rightarrow 0, mrd \rightarrow 0$  then  $H^*(\Delta_L, \epsilon, r) \cong H^*(\Delta_L)$

Using Lemma 2 for  $\overset{\wedge}{\Delta}_{L-\epsilon}(m, r)$  and  $\overset{\wedge}{\Delta}_{L+\epsilon}(m, r)$ , we have:

$$\text{IP}(\overset{\wedge}{\Delta}_{L+\frac{3}{2}\epsilon}(2r)) \subset \overset{\wedge}{\Delta}_{L+\epsilon}(m, r) \subset \overset{\wedge}{\Delta}_{L+\frac{1}{2}\epsilon}(2r) \geq 1 - 3m e^{-\frac{C}{8r^2} mrd}$$

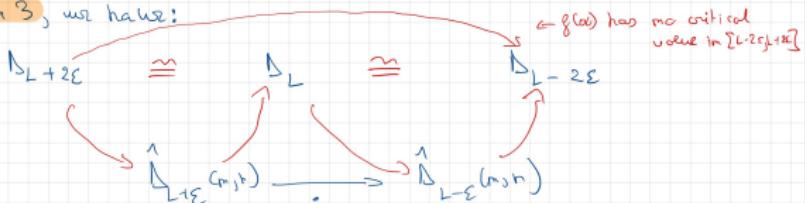
$$\text{IP}(\overset{\wedge}{\Delta}_{L-\frac{1}{2}\epsilon}(2r)) \subset \overset{\wedge}{\Delta}_{L-\epsilon}(m, r) \subset \overset{\wedge}{\Delta}_{L-\frac{3}{2}\epsilon}(2r) \geq 1 - 3m e^{-\frac{C}{8r^2} mrd}$$

We want  $\overset{\wedge}{\Delta}_L$  not  $\overset{\wedge}{\Delta}_L^{\text{ab}}$  but since we assume  $L$  is  $\epsilon$ -regular, if  $r$  small enough:

$$\overset{\wedge}{\Delta}_{L+2\epsilon} \subset \overset{\wedge}{\Delta}_{L+\frac{3}{2}\epsilon} \subset \overset{\wedge}{\Delta}_{L+\frac{1}{2}\epsilon}(2r) \subset \overset{\wedge}{\Delta}_L \subset \overset{\wedge}{\Delta}_{L-\frac{1}{2}\epsilon} \subset \overset{\wedge}{\Delta}_{L-\frac{3}{2}\epsilon}(2r) \subset \overset{\wedge}{\Delta}_{L-2\epsilon}$$

$$\text{So we have: } \text{IP}(\overset{\wedge}{\Delta}_{L+2\epsilon} \subset \overset{\wedge}{\Delta}_{L+\epsilon}(m, r) \subset \overset{\wedge}{\Delta}_L \subset \overset{\wedge}{\Delta}_{L-\epsilon}(m, r) \subset \overset{\wedge}{\Delta}_{L-2\epsilon}) \geq 1 - 6m e^{-\frac{C}{8r^2} mrd}$$

Using Lemma 3, we have:



by Lemma 3:  $H_*(\overset{\wedge}{\Delta}_L) \cong H_*(\text{Im}(i_*))$

# Application to manifold learning



Sorbonne University

Let  $\mathcal{M}$  be a smooth  $m$ -dimensional closed manifold in  $\mathbb{R}^d$ . We wish to recover the homology of  $\mathcal{M}$  given  $\mathcal{X}_n = \{X_1, \dots, X_n\}$ .

→ The case  $\mathcal{X}_n \subset \mathcal{M}$  have been extensively studied [ref], for example:

→ **Theorem** : If  $nr^d \geq C \log(n)$  and  $C > (w_d p_{min})^{-1}$  then  
 $H_*(U(\mathcal{X}_n, r)) \cong H_*(\mathcal{M})$  as.

→ Is there similar result when noise is present ?

→ Yes in special case or in the course with wasserstein.

**Goal :** We want to resolve this problem with a large class of distributions and with as few assumptions as possible.

# Application to manifold learning



Sorbonne University

Let  $\mathcal{M}$  be a smooth  $m$ -dimensional closed manifold in  $\mathbb{R}^d$ . We wish to recover the homology of  $\mathcal{M}$  given  $\mathcal{X}_n = \{X_1, \dots, X_n\}$ .

→ The case  $\mathcal{X}_n \subset \mathcal{M}$  have been extensively studied [ref], for example:

→ **Theorem** : If  $nr^d \geq C \log(n)$  and  $C > (w_d p_{min})^{-1}$  then  
 $H_*(U(\mathcal{X}_n, r)) \cong H_*(\mathcal{M})$  as.

→ Is there similar result when noise is present ?

→ Yes in special case or in the course with wasserstein.

**Goal :** We want to resolve this problem with a large class of distributions and with as few assumptions as possible.

# Application to manifold learning



Sorbonne University

Let  $\mathcal{M}$  be a smooth  $m$ -dimensional closed manifold in  $\mathbb{R}^d$ . We wish to recover the homology of  $\mathcal{M}$  given  $\mathcal{X}_n = \{X_1, \dots, X_n\}$ .

→ The case  $\mathcal{X}_n \subset \mathcal{M}$  have been extensively studied [ref], for example:

→ **Theorem** : If  $nr^d \geq C \log(n)$  and  $C > (w_d p_{min})^{-1}$  then  
 $H_*(U(\mathcal{X}_n, r)) \cong H_*(\mathcal{M})$  as.

→ Is there similar result when noise is present ?

→ Yes in special case or in the course with wasserstein.

**Goal :** We want to resolve this problem with a large class of distributions and with as few assumptions as possible.

# Application to manifold learning



Sorbonne University

Let  $\mathcal{M}$  be a smooth  $m$ -dimensional closed manifold in  $\mathbb{R}^d$ . We wish to recover the homology of  $\mathcal{M}$  given  $\mathcal{X}_n = \{X_1, \dots, X_n\}$ .

→ The case  $\mathcal{X}_n \subset \mathcal{M}$  have been extensively studied [ref], for example:

→ **Theorem** : If  $nr^d \geq C \log(n)$  and  $C > (w_d p_{min})^{-1}$  then  
 $H_*(U(\mathcal{X}_n, r)) \cong H_*(\mathcal{M})$  as.

→ Is there similar result when noise is present ?

→ Yes in special case or in the course with wasserstein.

**Goal :** We want to resolve this problem with a large class of distributions and with as few assumptions as possible.

# The general class of distribution

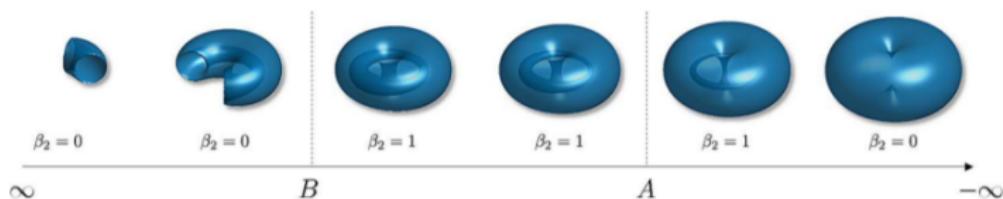


Sorbonne University

Let  $p : \mathbb{R}^d \rightarrow \mathbb{R}^+$  be a probability distribution. We say that  $p$  represents a noisy version of  $\mathcal{M}$ , if there exists  $0 < A < B < \infty$  st :

- $\forall L \in [A, B], D_L \cong \mathcal{M}$
- $\forall L > B, D_L \cong \mathcal{M}'$  where  $\mathcal{M}' \subset \mathcal{M}$  is a compact locally contractible proper subset of  $\mathcal{M}$

Exemple:



# The ideal case



Sorbonne University

If we know  $A, B$  then the recovery will be simple.

→ Choose  $L$  st  $[L - 2\epsilon, L + 2\epsilon] \subset [A, B]$

→ **Theorem 1** implies that  $H_*(L, \epsilon, n) \cong H_*(D_L) \cong H_*(M)$

We don't known  $A, B$ .

Solution :

- $M$  is a connected and orientable.
- $B - A > 8\epsilon$

# The ideal case



Sorbonne University

If we know  $A, B$  then the recovery will be simple.

→ Choose  $L$  st  $[L - 2\epsilon, L + 2\epsilon] \subset [A, B]$

→ Theorem 1 implies that  $H_*(L, \epsilon, n) \cong H_*(D_L) \cong H_*(M)$

We don't known  $A, B$ .

Solution :

- $M$  is a connected and orientable.
- $B - A > 8\epsilon$

# The procedure



Define:  $N_\epsilon = \sup [f(x)/2\epsilon]$        $L_{max} = 2\epsilon N_\epsilon$        $L_i = L_{max} - 2i\epsilon$

- 1** Compute  $H_*(L_i, \epsilon, n)$  for all  $i = 1, \dots, N_\epsilon$
- 2** Define  $i_* = 1 + \min\{i \in \{1, \dots, N_\epsilon\} : \text{rank}(H_*(L_i, \epsilon, n)) = 1\}$
- 3** Return  $H_*(L_{i_*}, \epsilon, n)$

**Remarks :** we have to choose r for the procedure, the next theorem tell us that if we choose r correctly, we have consistence.

# The procedure



Define:  $N_\epsilon = \sup [f(x)/2\epsilon]$        $L_{max} = 2\epsilon N_\epsilon$        $L_i = L_{max} - 2i\epsilon$

- 1** Compute  $H_*(L_i, \epsilon, n)$  for all  $i = 1, \dots, N_\epsilon$
- 2** Define  $i_* = 1 + \min\{i \in \{1, \dots, N_\epsilon\} : \text{rank}(H_*(L_i, \epsilon, n)) = 1\}$
- 3** Return  $H_*(L_{i_*}, \epsilon, n)$

**Remarks :** we have to choose  $r$  for the procedure, the next theorem tell us that if we choose  $r$  correctly, we have consistence.

# Theorem 2



Sorbonne University

## Theorem

Let  $\mathcal{M}$  be a  $m$ -dimensional closed, connected, orientable manifold in  $\mathbb{R}^d$ . Let  $X_1, \dots, X_n$  be sample from noisy density of  $\mathcal{M}$ . if we choose  $r \rightarrow 0$  such that  $nr^d \geq D \log(n)$  with  $D > (C_{\epsilon/2}^*)^{-1}$ . Applying the procedure then:

$$\mathbb{P}(H_*(L_{i*}, \epsilon, n) \cong H_*(\mathcal{M})) = 1$$

## Remarks:

- Need to know  $m, \epsilon, L_{max}, N_\epsilon$ .

# The proof of the theorem 2



Sorbonne University

Proof.

Recall that  $N_\varepsilon = \{f(x)/2\varepsilon\}$  and  $L_{\max} = 2\varepsilon N_\varepsilon$ .

$$\text{let } E = \bigcup_{i=1}^{N_\varepsilon} L_{l_{i-1}} = L_{l_1} + 2\varepsilon \hookrightarrow \bigcup_{i=1}^{\lceil m/r \rceil} L_{l_i} + \varepsilon \hookrightarrow \bigcup_{i=1}^{\lceil m/r \rceil} L_{l_i} \hookrightarrow \bigcup_{i=1}^{\lceil m/r \rceil} L_{l_i} - \varepsilon = \bigcup_{i=1}^{\lceil m/r \rceil} L_{l_{i-1}}$$

Applying Lemma 2

$$P(E) \geq 1 - 3N_\varepsilon m^{-\frac{C}{\varepsilon} m + d}$$

Now, we have to show that  $[L_{l_1} - 2\varepsilon, L_{l_1} + 2\varepsilon] \subset (A, B)$  Pcaré

Rq: Since we assume that  $M$  is connected we have  $\beta_m(M) = 1 = \beta_m(M')$

and if  $L > B$   $\beta_m(M') = 0$

Our requirement that  $L_{l_1} - L_{l_0} = 2\varepsilon$  and  $B - A \geq 8\varepsilon$  guarantees that there are at least four consecutive  $L_{l_1} > L_{l_2} > L_{l_3} > L_{l_4} \in (A, B)$ .

For  $k=2, 3$  we have  $[L_{l_k} - 2\varepsilon, L_{l_k} + 2\varepsilon] \subset (A, B)$  but it is not true for  $i_1$ .

Since  $\hat{\beta}_m(L_{l_1}, \varepsilon, m) = 1 \text{ or } 0$  then  $i_+$  is  $i_2$  or  $i_3$  but in both cases we have

$$[L_{l_1} - 2\varepsilon, L_{l_1} + 2\varepsilon] \subset (A, B)$$

$PH_*(f)$  : tracks the evolution of the homology of  $D_L$  for  $L = +\infty$  to  $-\infty$   
→ We have shown that we can recover  $D_L$  for all  $L > 0$   
→ But  $L$  is continuous

However, we propose this procedure : Let  $N_\epsilon, L_{max}, L_i$  like previously  
Define  $\hat{D}_\epsilon = \{\hat{D}_{L_i}(n, r)\}_{i \in \mathbb{Z}}$  and  $\hat{PH}^\epsilon(f)$  the persistent homology of  $\hat{D}_\epsilon$  as  
the estimator of  $PH_*(f)$

## Theorem

if  $r \rightarrow 0$  and  $nr^d \rightarrow \infty$  then :

$$\mathbb{P}(d_B(\hat{PH}^\epsilon(f), PH_*(f)) \leq 5\epsilon) \geq 1 - 3N_\epsilon n e^{-C_{\epsilon/2} nr^d}$$

In particular, if ....

How we compute  $H_*(U(\mathcal{X}, r))$  ?

→ nerf theorem  $\Rightarrow H_*(\text{Cech}(\mathcal{X}), r) \cong H_*(U(\mathcal{X}), r)$   $\Rightarrow$  Theorems holds for  $\text{Cech}(\mathcal{X}), r$

→ But  $H_*(\text{Cech}(\mathcal{X}), r)$  is too cost computationally

Like in the course, we want to use an approximation of the Cech: Rips.  
But nerf theorem don't holds for Rips, however we can show that the theorems still holds for Rips.

How we compute  $H_*(U(\mathcal{X}, r))$  ?

→ nerf theorem  $\Rightarrow H_*(\text{Cech}(\mathcal{X}), r) \cong H_*(U(\mathcal{X}), r)$   $\Rightarrow$  Theorems holds for  $\text{Cech}(\mathcal{X}), r$ )

→ But  $H_*(\text{Cech}(\mathcal{X}), r)$  is too cost computationally

Like in the course, we want to use an approximation of the Cech: Rips.  
But nerf theorem don't holds for Rips, however we can show that the theorems still holds for Rips.

How we compute  $H_*(U(\mathcal{X}, r))$  ?

→ nerf theorem  $\Rightarrow H_*(\text{Cech}(\mathcal{X}), r) \cong H_*(U(\mathcal{X}), r)$   $\Rightarrow$  Theorems holds for  $\text{Cech}(\mathcal{X}), r$ )

→ But  $H_*(\text{Cech}(\mathcal{X}), r)$  is too cost computationally

Like in the course, we want to use an approximation of the Cech: Rips.  
But nerf theorem don't holds for Rips, however we can show that the theorems still holds for Rips.

How we compute  $H_*(U(\mathcal{X}, r))$  ?

→ nerf theorem  $\Rightarrow H_*(\text{Cech}(\mathcal{X}), r) \cong H_*(U(\mathcal{X}), r)$   $\Rightarrow$  Theorems holds for  $\text{Cech}(\mathcal{X}), r$

→ But  $H_*(\text{Cech}(\mathcal{X}), r)$  is too cost computationally

Like in the course, we want to use an approximation of the Cech: Rips.  
But nerf theorem don't holds for Rips, however we can show that the theorems still holds for Rips.

# Compute the homology of the image map



Sorbonne University

First let compute  $\beta_k(\Delta)$  for  $\Delta$  a simplicial complex.

- $\Delta_k = \{\sigma_1, \dots, \sigma_2\}$
- The map  $T_k : \Delta_k^\pi \rightarrow \mathbb{R}^k \quad \forall \sigma_i \in \Delta_k \quad T_k(\sigma_i) = \text{sgn}(\pi)e_i$

We can define the boundary matrix  $\partial_k$  be  $n_{k-1} \times n_k$  st :

$$(\partial_k)_i = \sum_{\sigma \in \Delta_{k-1}} T_k(\sigma)(\text{image})$$

So  $L_k = \partial_{k+1}\partial_{k+1}^T + \partial_k\partial_k^T$  Then  $\beta_k(\Delta) = \text{rank}(\ker(L_k))$

Example:

$$\text{Let } C_n = \{\sigma_1, \sigma_2, \dots, \sigma_{m_n}\} \text{ and } C_{k-1} = \{\tau_1, \tau_2, \dots, \tau_{m_{k-1}}\}$$

$$\partial_k = \begin{pmatrix} \tau_1 & \tau_2 & \cdots & \tau_{m_{k-1}} \\ a_{11}^1 & a_{12}^1 & \cdots & a_{1m_{k-1}}^1 \\ a_{21}^1 & a_{22}^1 & \cdots & a_{2m_{k-1}}^1 \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \tau_{m_{k-1}} & a_{m_{k-1}1}^2 & a_{m_{k-1}2}^2 & \cdots & a_{m_{k-1}m_{k-1}}^2 \end{pmatrix}$$

Where  $a_{ij}^k$  if the  $i$ -th complete of  $\tau_{k-1}$  is a face of  $j$ -complete of  $C_k$ .



# Compute the homology of the image map



Sorbonne University

Let  $\Delta^1$ ,  $\Delta^2$  and the map  $i_k : H_k(\Delta^1) \rightarrow H_k(\Delta^2)$ , our goal is to compute the homology of  $Im(i_k)$ .

- $\Delta_k^1 = \{\sigma_1, \dots, \sigma_{n_k^1}\}$
- $\Delta_k^2 = \{\sigma_1, \dots, \sigma_{n_k}, \sigma_{n_k+1}, \dots, \sigma_{n_k^2}\}$  -

$$\partial_k^{(2)} = \begin{bmatrix} \partial_k^{(1)} & \cdots \\ 0 & \ddots \end{bmatrix}$$

Now, if  $\{v_1, \dots, v_m\} \subset \mathbb{R}^{n_k^1}$  is a basis for  $\ker(L_k^1)$ . Let  $\hat{v}_i \in \mathbb{R}^{n_k^2}$  be zeros padded version of  $v_i$  by linear algebra, we can find that:

$$rank(Im(i_k)) = rank(\hat{\partial}_{k+1}^{(2)}) - rank(\partial_{k+1}^{(2)})$$

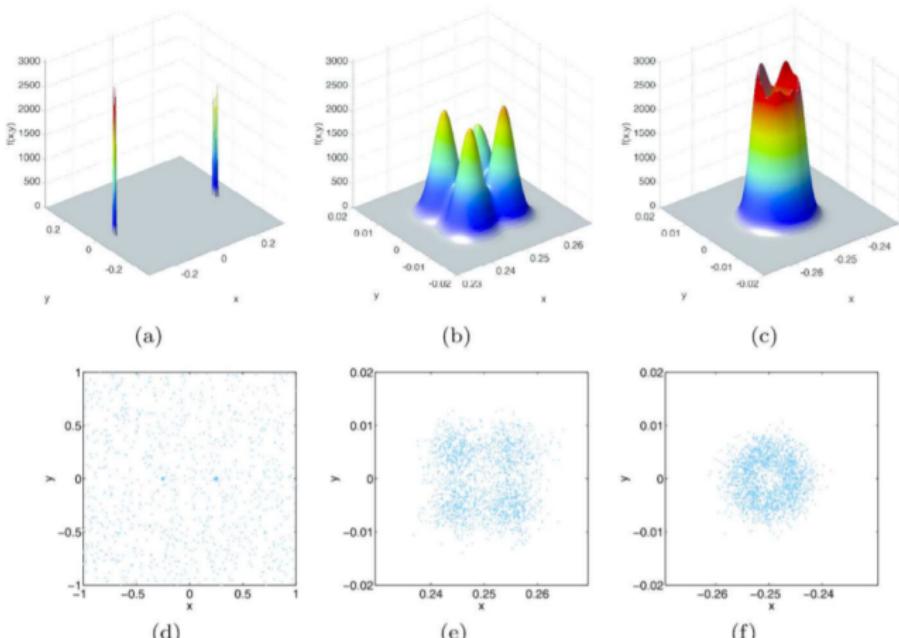
where  $\hat{\partial}_{k+1}^{(2)} = (\partial_{k+1}^{(2)}, \hat{v}_1, \dots, \hat{v}_m)$

# Results on simulated data



Sorbonne University

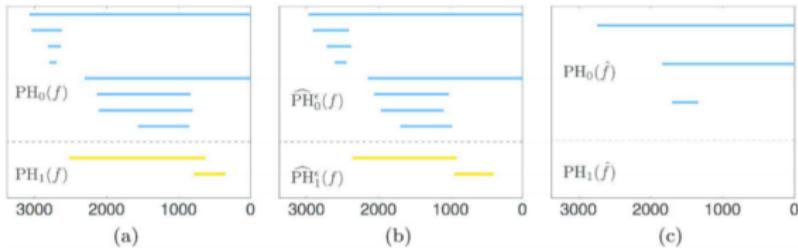
- In the papers, there are results on binary classification, kernel regression and spectral clustering.
- We decide to show an application for Hierarchical clustering :



# Results of the estimators



Sorbonne University



**Figure 14.** Estimating the persistent homology of the density function  $f$  presented in Figure 13.  
(a) The “true” barcode for the function  $f$ , that is,  $\text{PH}_*(f)$  (computed by sampling the density function on a fine grid). (b) The barcode computed from the estimator  $\widehat{\text{PH}}_*^\varepsilon(f)$ . The parameters used are  $n = 5000$ ,  $r = 0.001$ ,  $\varepsilon = 3.5$ . (c) The barcode computed for the kernel density estimator –  $\text{PH}_*(\hat{f})$ . The kernel parameters are the same as for  $\widehat{\text{PH}}_*^\varepsilon(f)$ , the grid size taken is  $500 \times 500$ . Note that the estimator  $\widehat{\text{PH}}_*^\varepsilon(f)$  gives a result that is very similar to the true barcode. In both cases there are five significant features in  $H_0$  and two significant features in  $H_1$ . The barcode for  $\text{PH}_*(\hat{f})$  only recover the coarse features, namely the two clusters, but completely ignores the finer structures. We note that for visualization purposes we filtered out the very small bars before drawing the barcodes here.

# Conclusion



Sorbonne University

Thank you !

Source :

- Topological consistency via kernel estimation [Omer Bobrowski, Sayan Mukherjee, Jonathan E. Taylor]