

TDA - fall 2025

bottleneck
distance +
stability



Recap

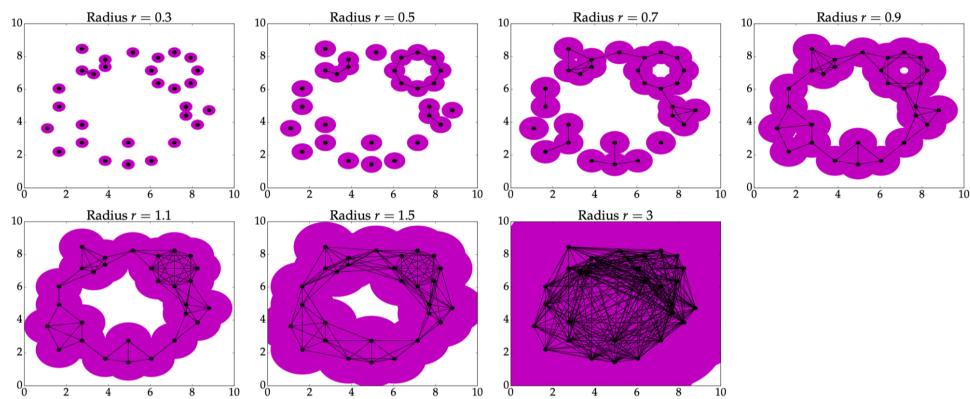
- Welcome back!
 - ↳ How did you like AARTN talks?
- Next HW: a "paper chase"
 - Intention:
 - exposure to topics you like within TDA
 - practice reviewing
 - explore topics for final project (later on)

The workflow so far

Build a filtration F from a simplicial complex
↳ usually parameterized by a function
 f , via sublevel sets

Example: $P \subseteq \mathbb{R}^n$

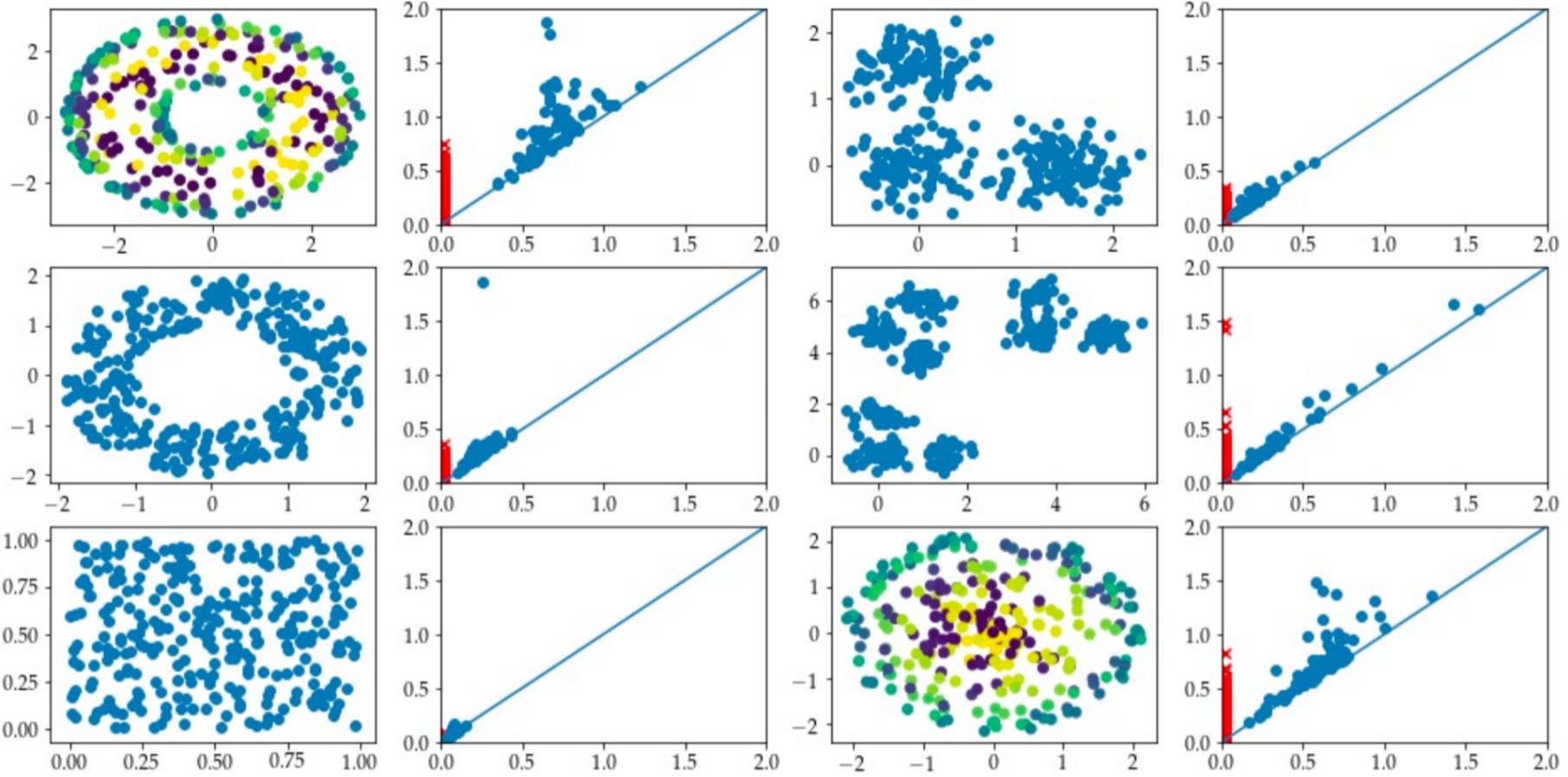
Build Rips filtration:
for $0 \leq r_1 \leq \dots \leq r_k$, $K_i = VR(R, r_i)$



→ Persistence diagram

Result!

H_0 and H_1



... now what?

Distance measures

A **distance** on a set X is a function

$$d: X \times X \rightarrow \mathbb{R}_{\geq 0} \text{ s.t. } \forall x, y, z \in X$$

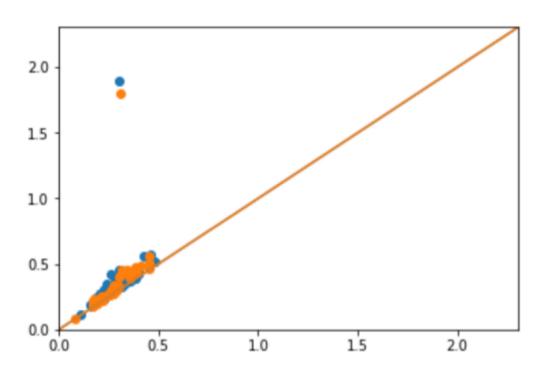
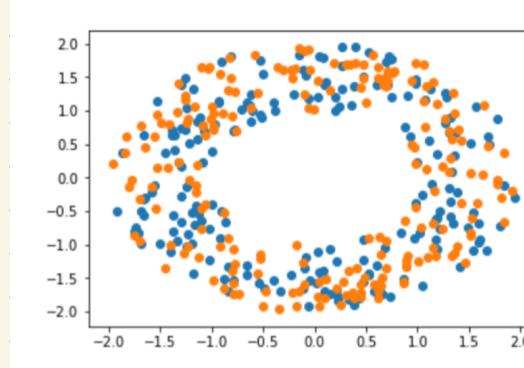
- $d(x, y) \geq 0$ + $d(x, y) = 0 \Leftrightarrow x = y$

- $d(x, y) = d(y, x)$

- $d(x, z) \leq d(x, y) + d(y, z)$

Our goal: distances for PDs

$$X = \left\{ D_{gm} F(K) \right\}_{F, K}$$



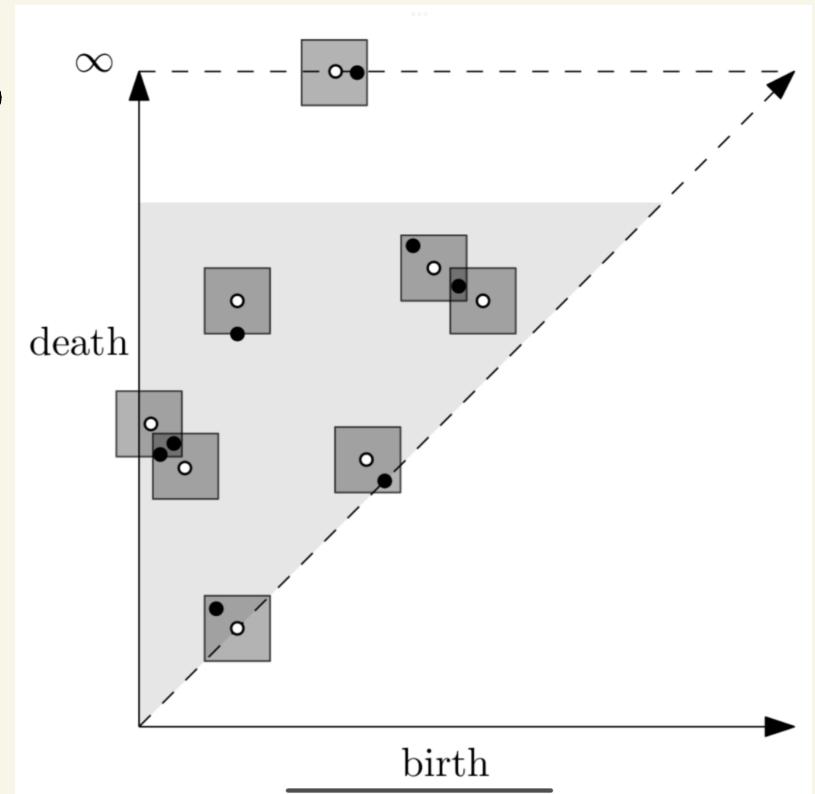
Bottleneck distance (book's version)

Let $\Pi = \{\pi : Dgm_p(f) \rightarrow Dgm_p(g)\}$

denote the set of all bijections from PD of f to PD of g .

Let $\|x - y\|_\infty = \max\{|x_1 - x_2|, |y_1 - y_2|\}$, where $\infty - \infty = 0$.

Then $d_B(Dgm(f), Dgm(g)) = \inf_{\pi \in \Pi} \sup_{x \in Dgm(f)} \|x - \pi(x)\|_\infty$



Fact: d_B is a metric.

Proof:

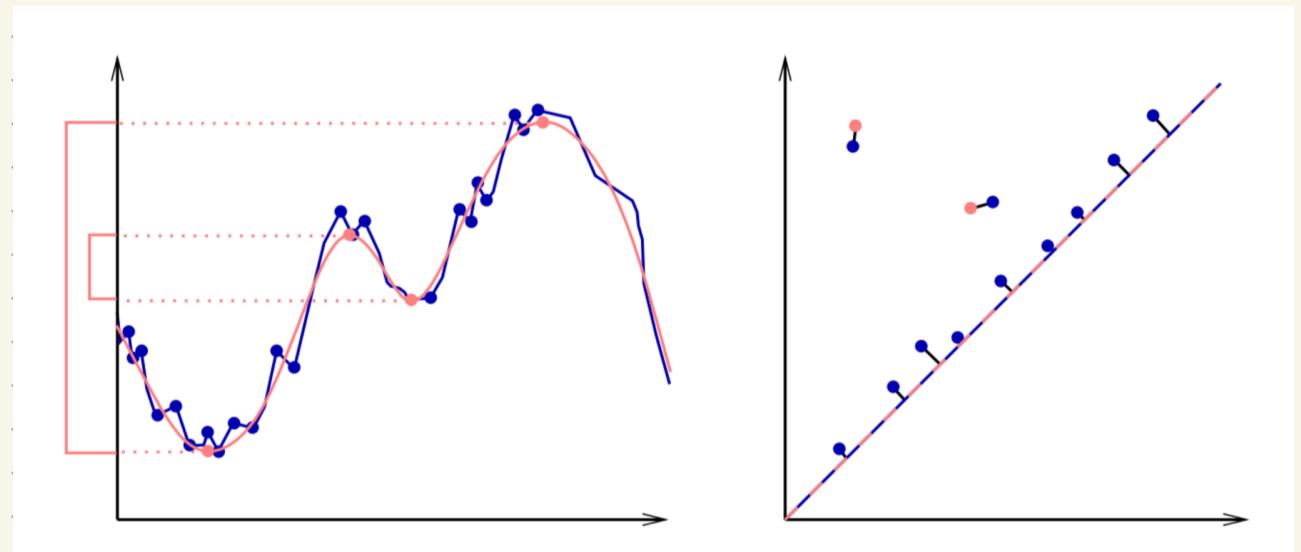
Stability

Let X be a triangulatable space &
 $f, g: X \rightarrow \mathbb{R}$ be tame functions
↳ (e.g.: smooth & Lipschitz)

giving rise to two space filtrations

F_f & F_g . Then $\forall p \geq 0$

$$d_B(Dgm(F_f), Dgm(F_g)) \leq \|f - g\|_\infty$$

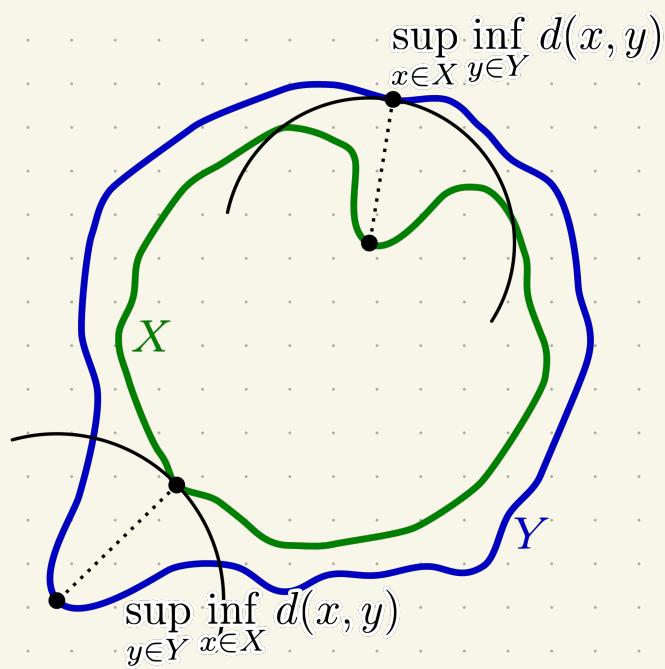


More stability

Let d_H be Hausdorff distance:

$$d_H(X, Y) = \max \left\{ \sup_x \inf_y \|x - y\|, \sup_y \inf_x \|x - y\| \right\}$$

Translating:



For finite point clouds $X, Y \subseteq \mathbb{R}^d$,
 let $Dgm(C(X))$ & $Dgm(C(Y))$ be the
 persistence diagrams of the filtration
 defined by the Čech complex. Then,
 $d_B(Dgm(C(X)), Dgm(C(Y)))$
 $\leq d_H(X, Y)$.

Proof picture:

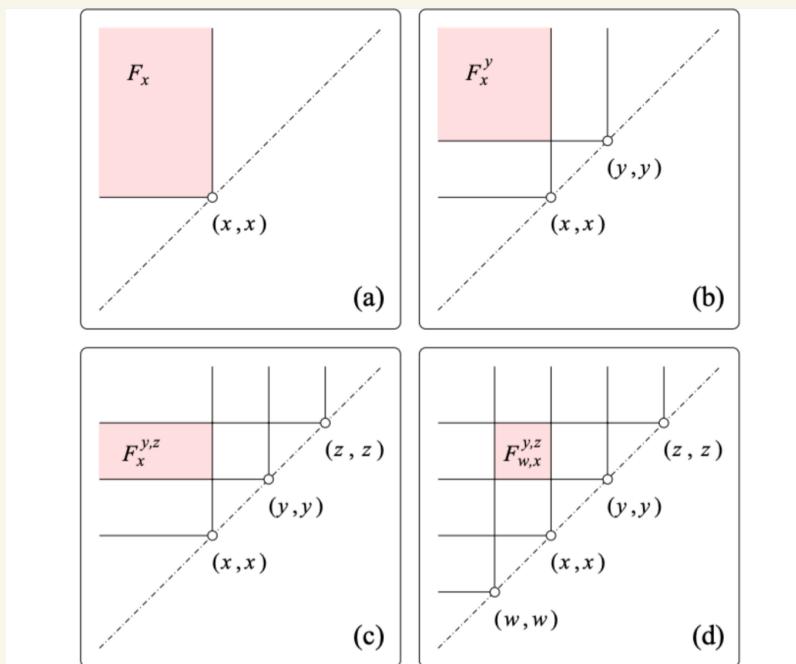


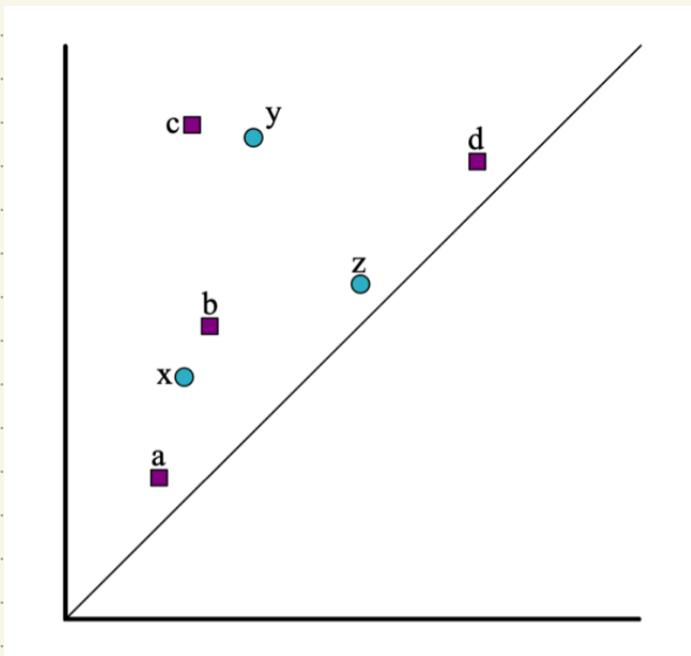
Figure 3: (a) Homology group of the sub-level set $f^{-1}(-\infty, x]$.
 (b) Image of F_x in F_y . (c) Kernel of surjection $F_x^y \rightarrow F_x^z$. (d)
 Quotient of $F_x^{y,z}$ and $F_w^{y,z}$.

P,q-Wasserstein Distance

Given diagrams $X + Y$,

$$W_p^q(X, Y) = \inf_{\ell: X \rightarrow Y} \left(\sum_{x \in X} \|x - \ell(x)\|_p^q \right)^{1/q}$$

$$= \inf_{\ell: X \rightarrow Y} \left(\sum_{x \in X} (\|x - \ell(x)\|_p^q)^{1/q} \right)^q$$



Special cases:

$$p = q = \infty :$$

$$p = q = 2 : \text{common!}$$

Note of warning:
Notation is not consistent!

Our book:

Definition 3.10 (Wasserstein distance). Let Π be the set of bijections as defined in Definition 3.9.
For any $p \geq 0, q \geq 1$, the q -Wasserstein distance is define as

$$d_{W,q}(\text{Dgm}_p(\mathcal{F}_f), \text{Dgm}_p(\mathcal{F}_g)) = \inf_{\pi \in \Pi} \left[\sum_{x \in \text{Dgm}_p(\mathcal{F}_f)} (\|x - \pi(x)\|_q)^q \right]^{1/q}.$$

Other reference:

$$W_q(X, Y) = \left[\inf_{\eta: X \rightarrow Y} \sum_{x \in X} \|x - \eta(x)\|_\infty^q \right]^{1/q}$$

Nice resource: AARSTN talk by
Kate Turner

Nonetheless, can get some weaker notion
of stability?

usually need addition of Lipschitz:

$$|f(x) - f(y)| \leq \|x - y\|$$

Then: $\exists C & k \geq 1$ s.t.

$$W_g^q(X, Y) \leq C \cdot \|f - g\|_{\infty}^{1 - \frac{k}{q}}$$

[Note: Hiding some technicalities here -

I recommend Skraba & Turner 2020

[if you are curious!]

Space of persistent diagrams

Let D_\emptyset be the empty diagram.

The space of persistent diagrams D_P^q
is the set of diagrams with
finite distance to D_\emptyset , ie

$$D_P^q = \{x \mid W_p^q(x, D_\emptyset) < \infty\}$$

so for each $x \in X$,

Note: does not necessarily mean
 X is finite!

Some statistical things

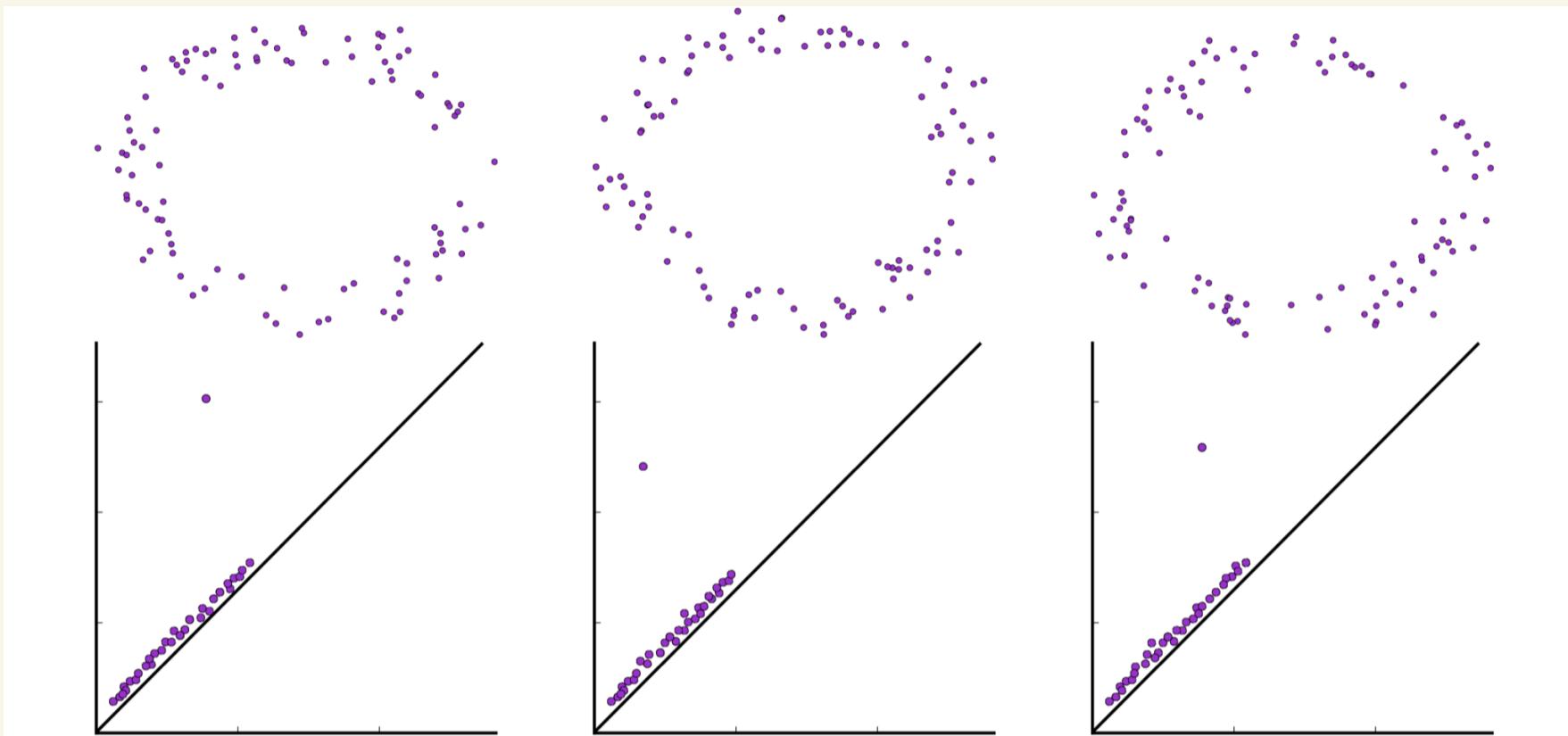
- D_p^q is complete & separable (Polish)
if $p = \infty$ & $q \in \mathbb{Z} \geq 1$

↳ Why?

- Non-negatively curved Alexandrov space if $p = q = 2$.

↳ Why?

How can we get an "average"?



Frechet means

Consider $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^2$

The Frechet variance of X is

$$\text{Var}_F = \inf_{a \in \mathbb{R}^2} \left\{ F_F(a) = \frac{1}{n} \sum_i \|x_i - a\|^2 \right\}$$

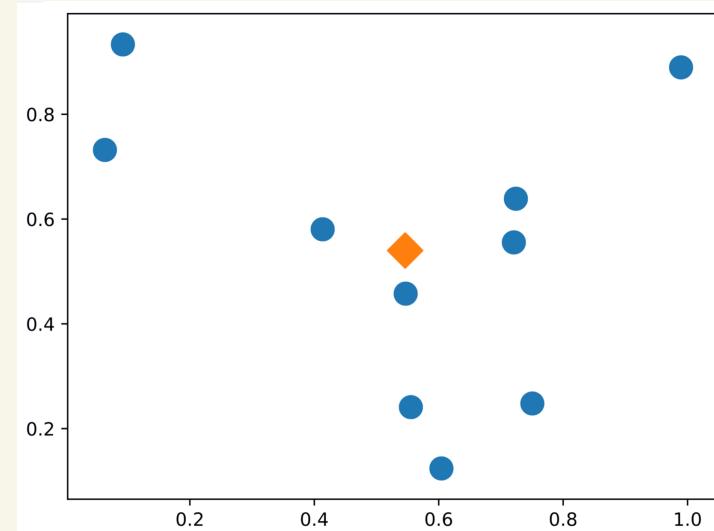
The set which realizes is

$$E(F) = \{a \in \mathbb{R}^2 \mid F_F(a) = \text{Var}_F\}$$

↳ called Frechet mean

(or Frechet expectation)

→ Unique & Computable!



Now let $\{X_1, \dots, X_n\}$ be persistence diagrams in D_p^P .

Frechet variance

$$\text{Var}_p = \inf_{Y \in Y} \left\{ F_p(Y) = \frac{1}{n} \sum_{i=1}^n W_p(X_i, Y)^2 \right\}$$

& Frechet mean is the set where value is obtained:

$$E(Y) = \{Y | F_p(Y) = \text{Var}_p\}$$

... What??

Picture: not unique!

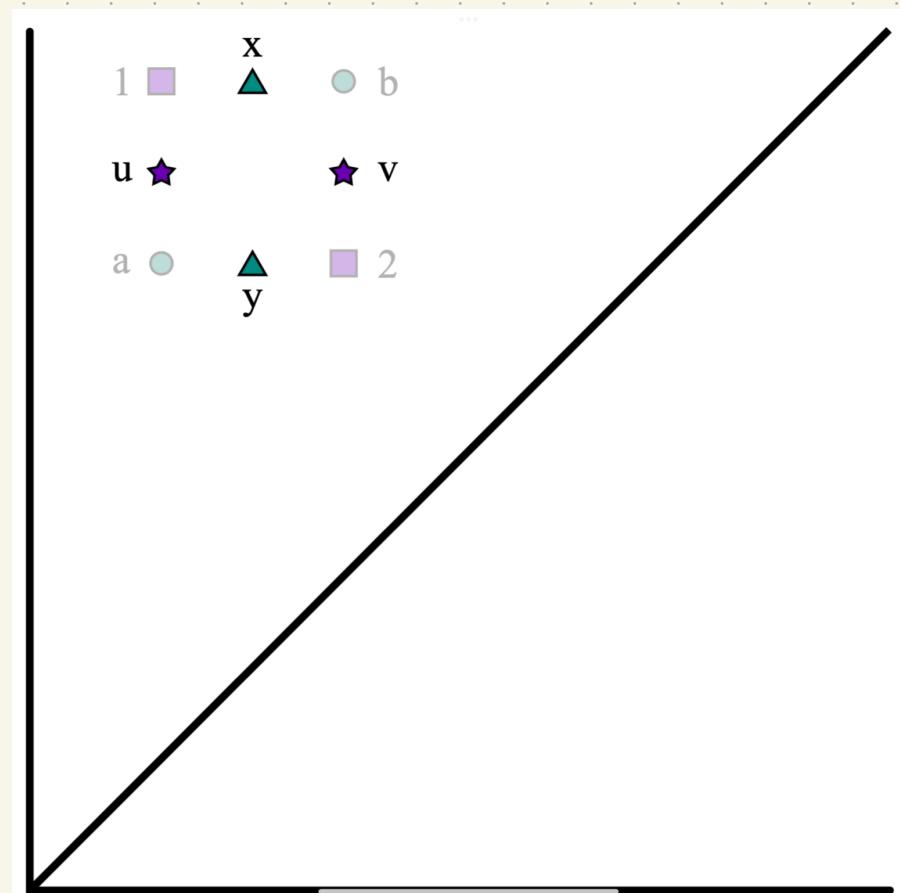
$$X_1 = \{\square_1, \square_2\}$$

$$X_2 = \{0_a, 0_b\}$$

Two Frechet means:

$$Y_1 =$$

$$Y_2 =$$



• $E(V) = \{Y_1, Y_2\}$

The good news:

The Frechet mean IS non-empty

Mileyko et al 2011
Turner et al 2014

(with some mild assumptions on
distribution of the set)

For D_2^2 : gradient descent algorithm
to compute local minimum.

In general, though:

Changing the question: Fasy et al 2014

What is an estimate for the average
(true?) diagram & how far off am I?

- Want to estimate PD

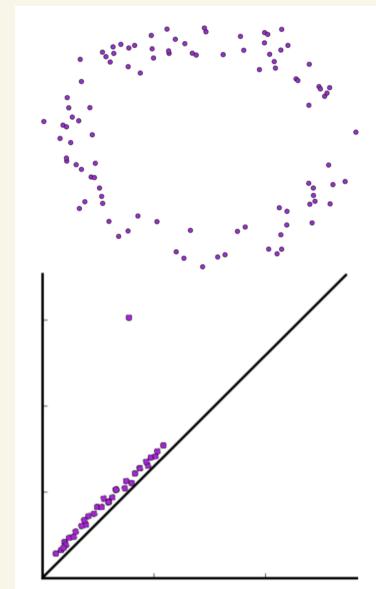
- for a set $M \subseteq \mathbb{R}^d$

- Don't know M

- But, have a sample

$S_n = \{x_1, \dots, x_n\}$ drawn uniformly from M .

- Persistence diagram for S_n is used as an estimator for $X \rightarrow$ denoted \hat{X}



Confidence Intervals

Given a collection of points $X = \{x_1, \dots, x_n\}$ from \mathbb{R} , the $100 \cdot (1-\alpha)\%$ confidence interval for the mean μ is the interval $[u(X), v(X)]$ such that

$$P(\mu \in [u(X), v(X)]) = 1 - \alpha$$

Equivalently: find c + an estimate for μ called $\hat{\mu}$ s.t.

$$P(|\mu - \hat{\mu}| \geq c) = \alpha$$

How to use in persistence?

Fix $\alpha \in (0, 1)$

Want $c_n := c_n(x_1, \dots, x_n)$ s.t.

$$\limsup_{n \rightarrow \infty} P(d_B(\hat{x}, x_n) > c_n) \leq \alpha$$

Then, $[0, c_n]$ is an asymptotic $(1-\alpha)$ confidence set for the bottleneck distance $d_B(\hat{x}, x)$.

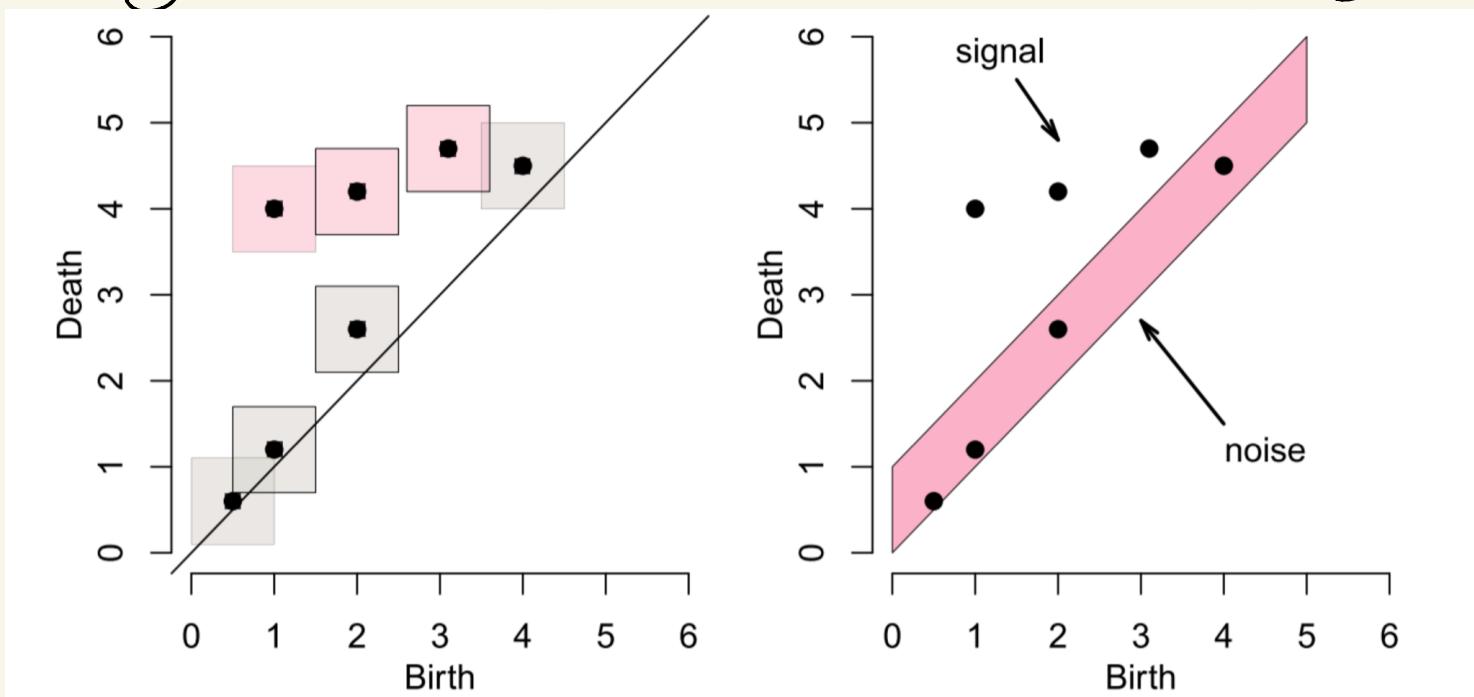
The confidence set C_n is the set of diagrams whose distance to \hat{x} is $\leq c_n$

$$C_n = \{Y \mid d_B(\hat{x}, Y) \leq c_n\}$$

Assume you have $\hat{X} + C_n$:

Put a box of width $2C_n$ at every point in \hat{X} .

A point is noise if its box intersects the diagonal \rightarrow or put strip along diagonal!



How to get C_b though?

- Start with data $S = \{x_1, \dots, x_n\}$
- Choose $b = b_n$ such that $b = O(\sqrt{n})$
- Pretend we have all $N = \binom{n}{b}$ subsamples S^1, \dots, S^N

↳ "bootstrapping"

(In reality: just do a lot)

- Calculate $d_+(S^j, S)$, $j = 1 \dots N$
- Set $L_b(t) = \frac{1}{N} \sum_{j=1}^N I(T_j > t)$
- Set $C_b = 2L_b^{-1}(\alpha)$

What now??

Using a theorem here!

Theorem

For mild assumptions on the space M , and for all large n ,

$$\mathbb{P}(d_B(\hat{X}, X) > c_b) \leq \mathbb{P}(d_H(S_n, M) > c_b) \leq \alpha + O\left(\frac{b}{n}\right)^{\frac{1}{4}}$$

[Note: there is every chance you
may be better at probability
than me.]

Back to pictures

(Next time—stay tuned!)