

TDA - fall 2025

Stats + ML
techniques
in persistence



Last time:

A Whole mess of statistics.

Prime theme:

Given a collection of PDFs, how

can we:

- Compute averages
- guarantee some amount of statistical significance
- eventually, maybe consider some methods that play well with ML

Changing the question: Fasy et al 2014

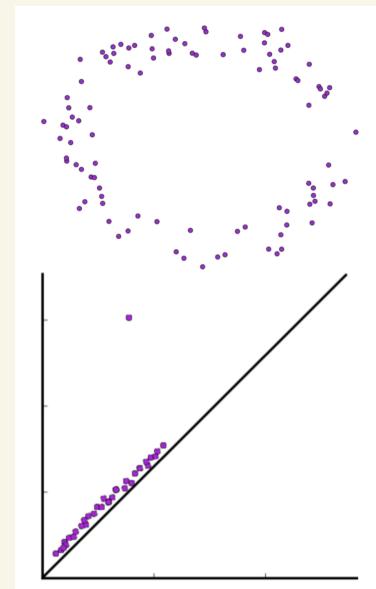
What is an estimate for the average
(true?) diagram & how far off am I?

- 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 $d_B(\hat{x}, x)$.

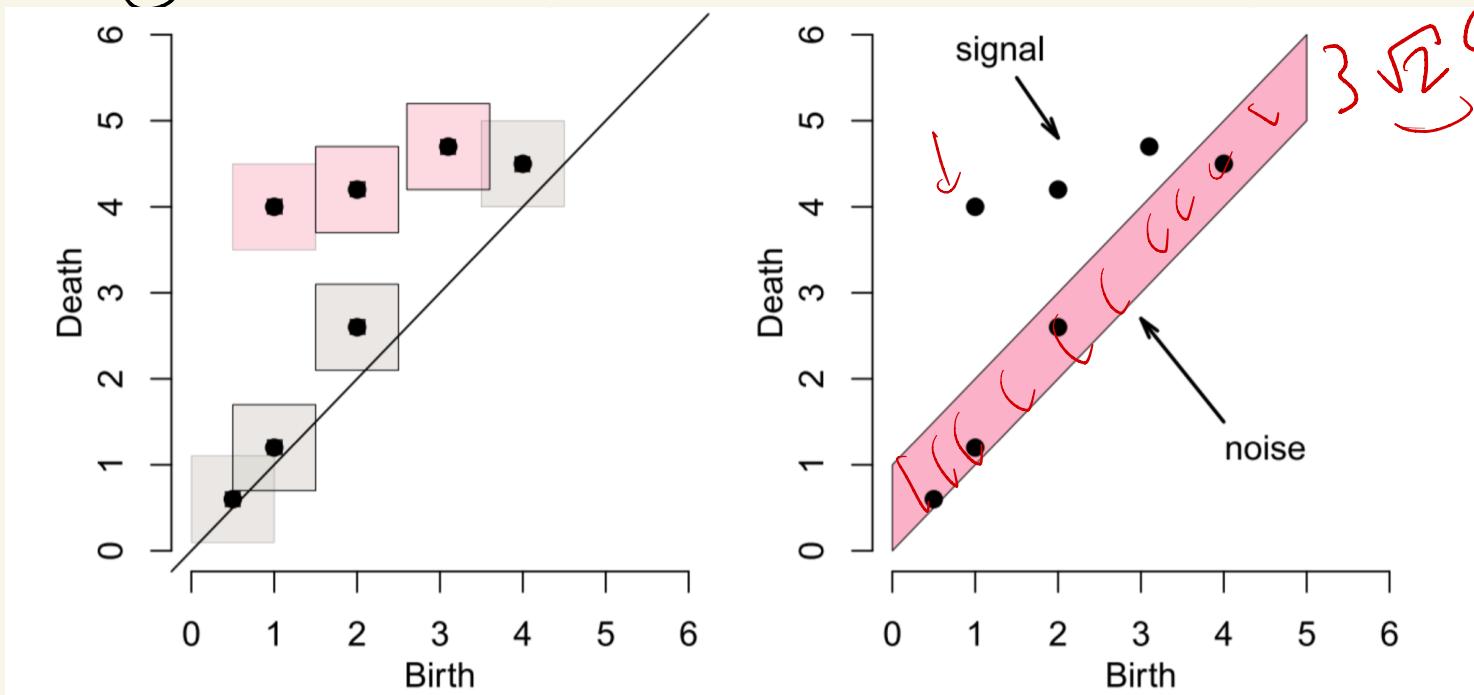
The confidence set C_n is the set of diagrams whose distance to \hat{x} is $\leq c_n$

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

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 few & pray)

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

$$\text{Set } C_b = 2L_b^{-1}(\alpha)$$

What now??

Using a theorem here!

Theorem

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

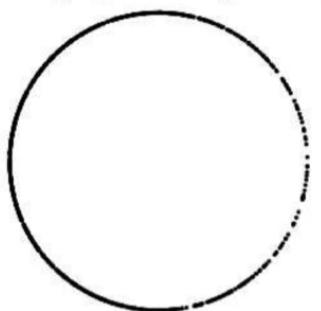
$$\mathbb{P}(d_B(\hat{X}, X) > c_b) \leq \mathbb{P}(d_H(S_n, M) > c_b) \leq \alpha + O\left(\frac{b}{n}\right)^{\frac{1}{4}}$$

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

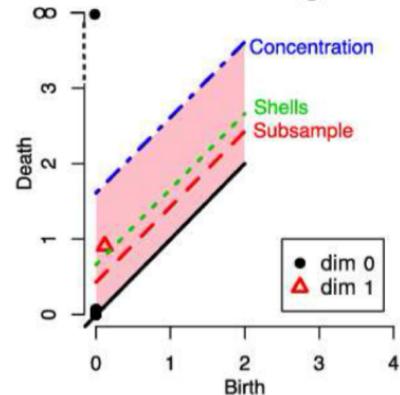
Some "toy examples" to demo:

Fasy et al
Annals of Statistics
2014

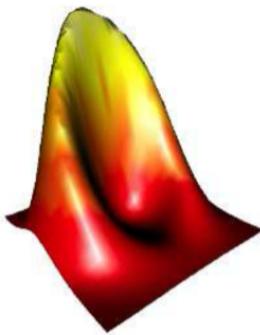
Circle ($r=1$) - Normal ($n=1000$)



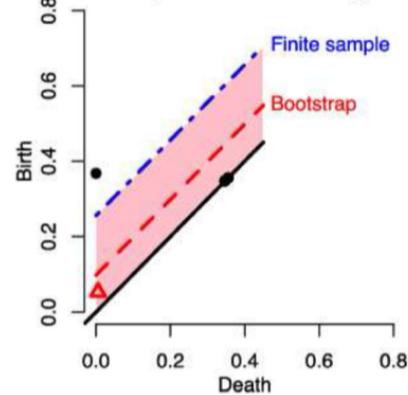
Persistence Diagram



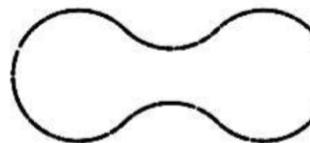
Kernel Density Estimator ($h=0.3$)



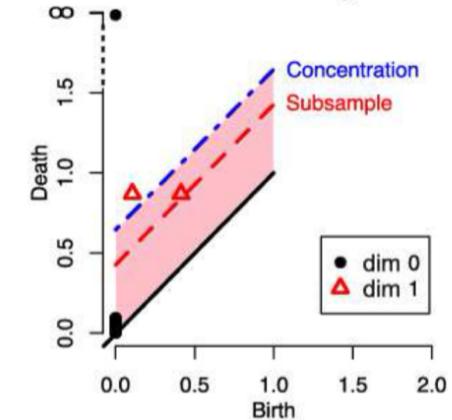
Density Persistence Diagram



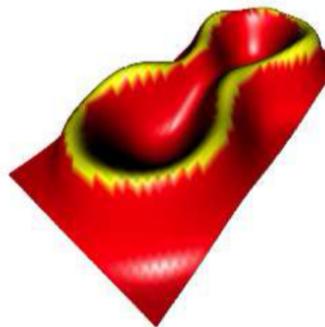
Eyeglasses - Uniform ($n=1000$)



Persistence Diagram



Kernel Density Estimator ($h=0.3$)



Density Persistence Diagram

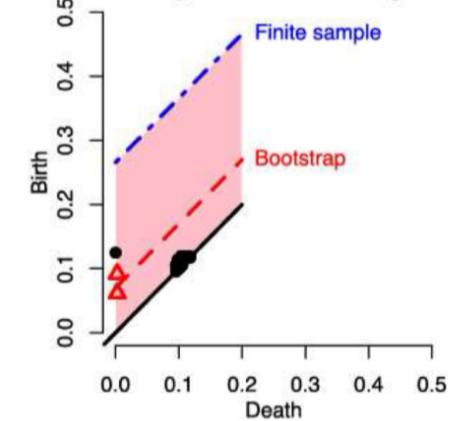


FIG. 7. Truncated Normal distribution over the unit Circle. (Top left) sample S_n . (Top right) corresponding persistence diagram. The black circles indicate the life span of connected components, and the red triangles indicate the life span of 1-dimensional holes. (Bottom left) kernel density estimator. (Bottom right) density persistence diagram. For more details see Example 14.

FIG. 8. Uniform distribution over the eyeglasses curve. (Top left) sample S_n . (Top right) corresponding persistence diagram. The black circles indicate the life span of connected components and the red triangles indicate the life span of 1-dimensional holes. Bottom left: kernel density estimator. (Bottom right) density persistence diagram. For more details see Example 15.

Some issues & takeaways

- Pros:
- Can in some sense prove what is noise versus a feature
 - Can provide some notion of average

- Cons:
- Averages are not unique
(even for simple diagrams)
 - Confidence intervals require diagrams which we get from subsamples of point clouds
↳ not always the data!
 - Understand how to use with ML
↳ what is an SVM for PDS?

Next: Changing the diagrams

We'll see 2 methods to alter diagrams so we can have better results:

- Persistence landscapes
- Persistence images

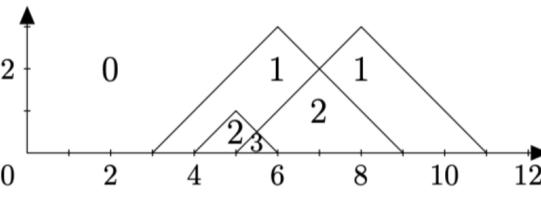
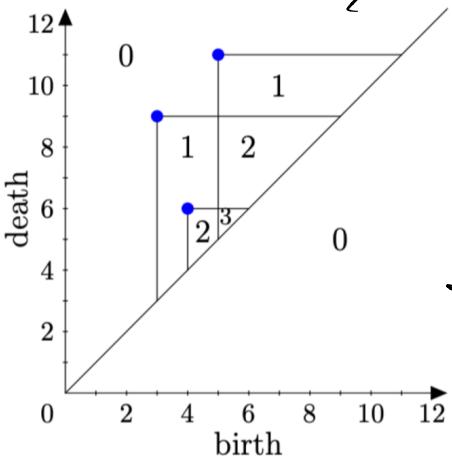
There are others! See Ch. 13 of book, or any NeurIPS or ICM2 paper mentioning "persistence" in the last 5 years. :)

Persistence Landscapes

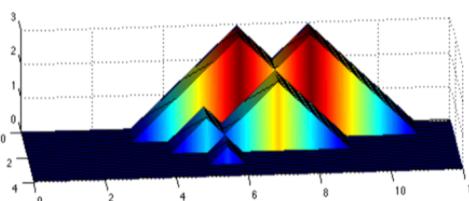
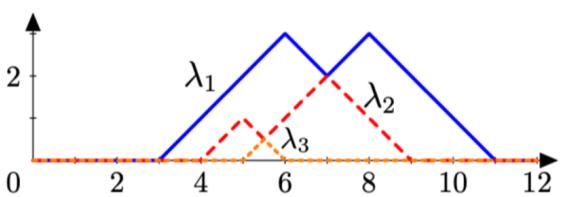
Bubenik, 2014 - JMLR

Definition (by picture): Fix dimension k .

step 3



Step 2



3) Rotate $(x, y) \mapsto \left(\frac{x+y}{2}, \frac{y-x}{2} \right)$,

$$\text{eg } (3, 9) \mapsto$$

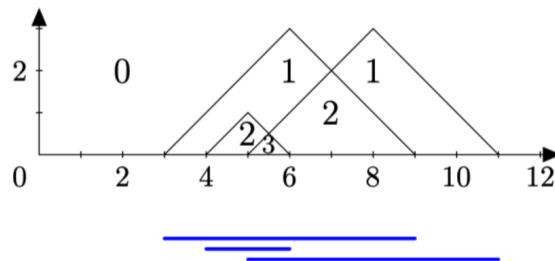
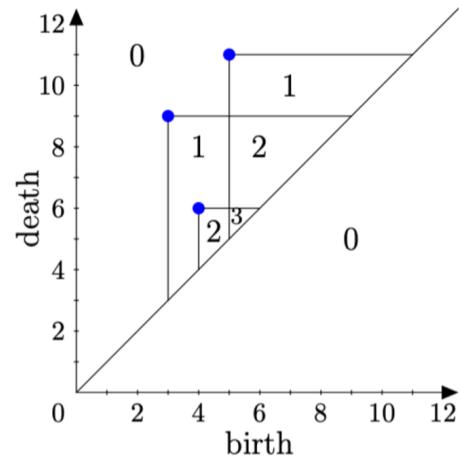
1) Compute the persistence diagram

2) Rank function:

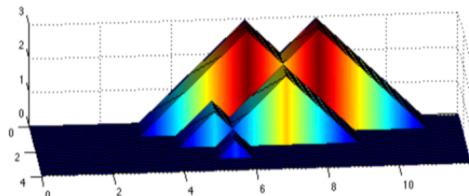
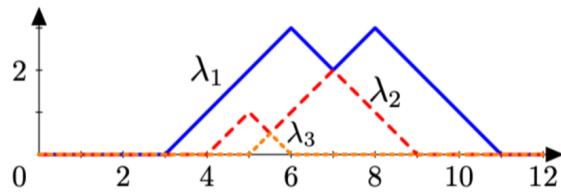
$$\beta^{a,b} = \dim(\text{Im}(H_k(X_a) \rightarrow H_k(X_b)))$$

(cont)

PLs (cont.)



So far:
here ✓



Step 4: Compute $\lambda: \mathbb{R}^2 \rightarrow \mathbb{R}$

$$\lambda(m, h) = \begin{cases} \beta_{m-h, m+h} & \text{if } h \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{i.e.: } \lambda(9, 1) =$$

[Calculate via barcodes if easier!]

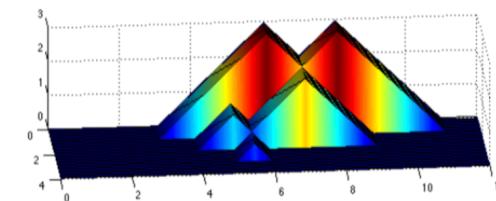
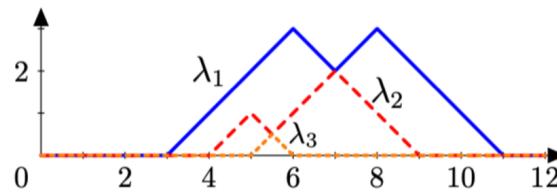
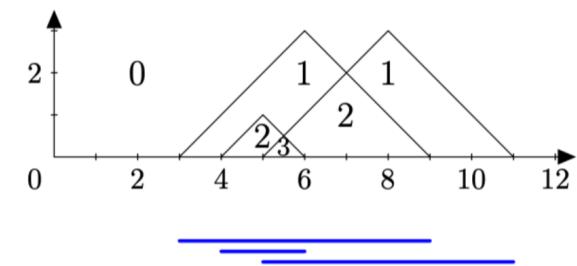
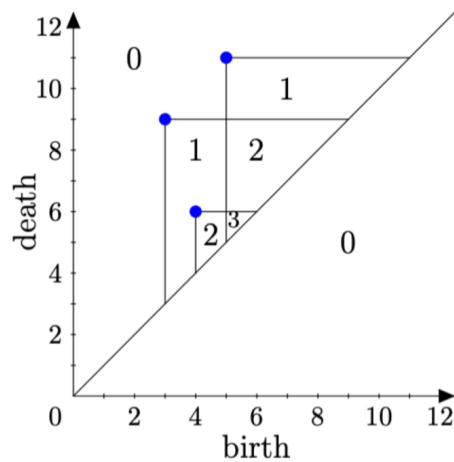
Step 5:

Next,

$$\text{let } \lambda_k(t) = \sup \left\{ m \geq 0 \mid \beta^{t-m, t+m} \geq k \right\}$$

so $\lambda_1 = \max$ s.t. Betti number is ≥ 1

$\lambda_2 = \max$ s.t. " " ≥ 2
etc.



Why?

Let's

Can get unique averages!
recall Fréchet means, & compare!

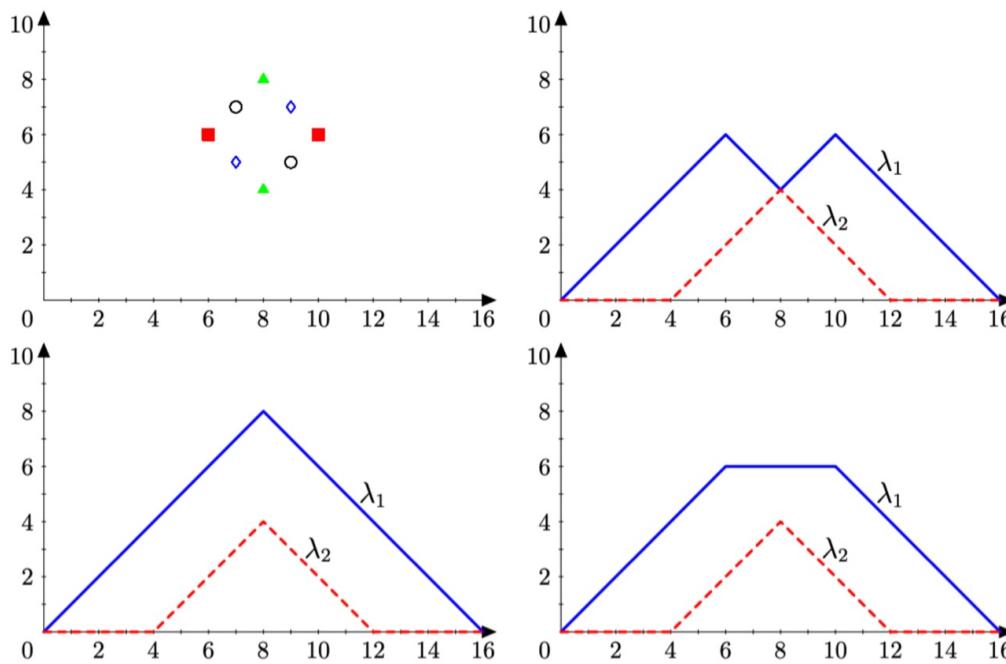
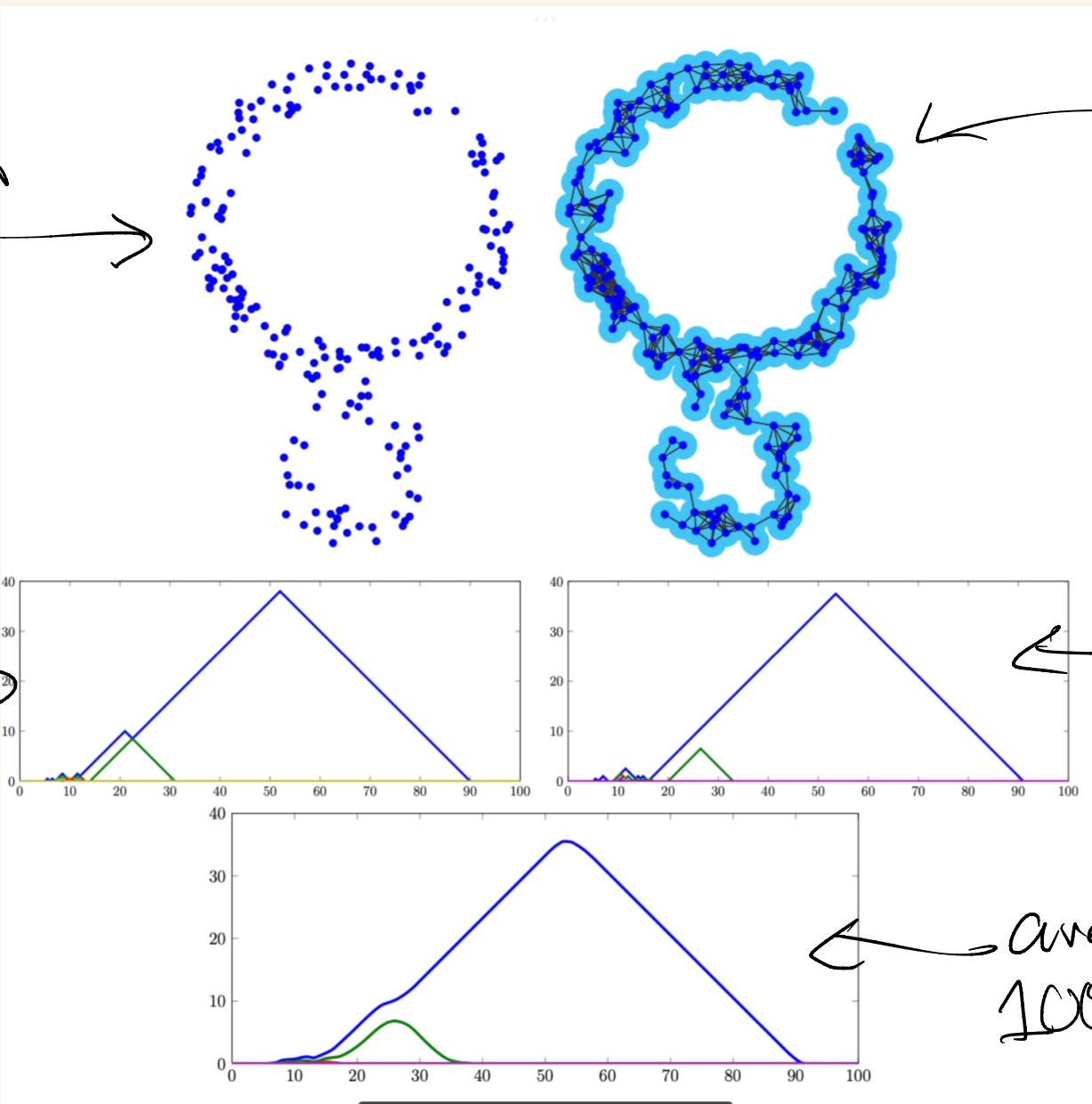


Figure 3: Means of persistence diagrams and persistence landscapes. Top left: the rescaled persistence diagrams $\{(6, 6), (10, 6)\}$ and $\{(8, 4), (8, 8)\}$ have two (Fréchet) means: $\{(7, 5), (9, 7)\}$ and $\{(7, 7), (9, 5)\}$. In contrast their corresponding persistence landscapes (top right and bottom left) have a unique mean (bottom right).

Can also derive (better) confidence intervals.

Toy examples : averages

Subsampled points from linked annuli



1 skeleton
of Čech
complex

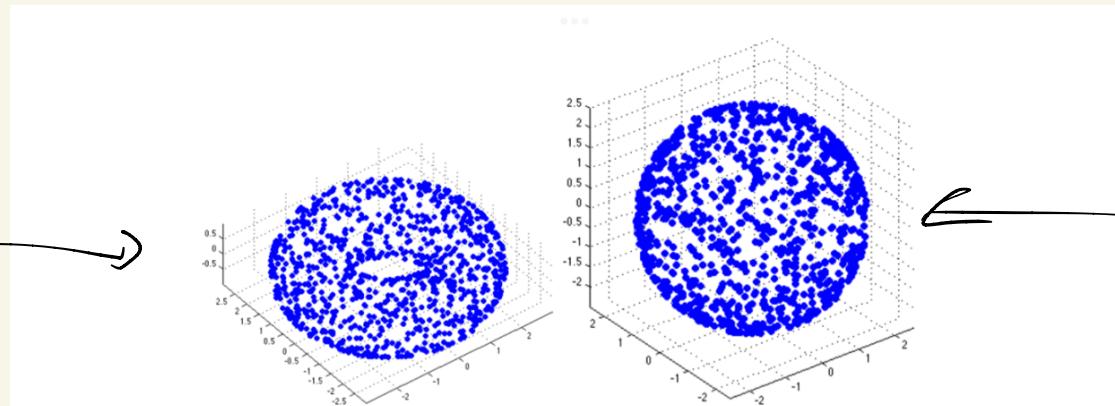
Sample 1
for H_1
 $P_{\bullet} L_0$

Sample 2

average of
100 samples

Another

Torus
Sampled



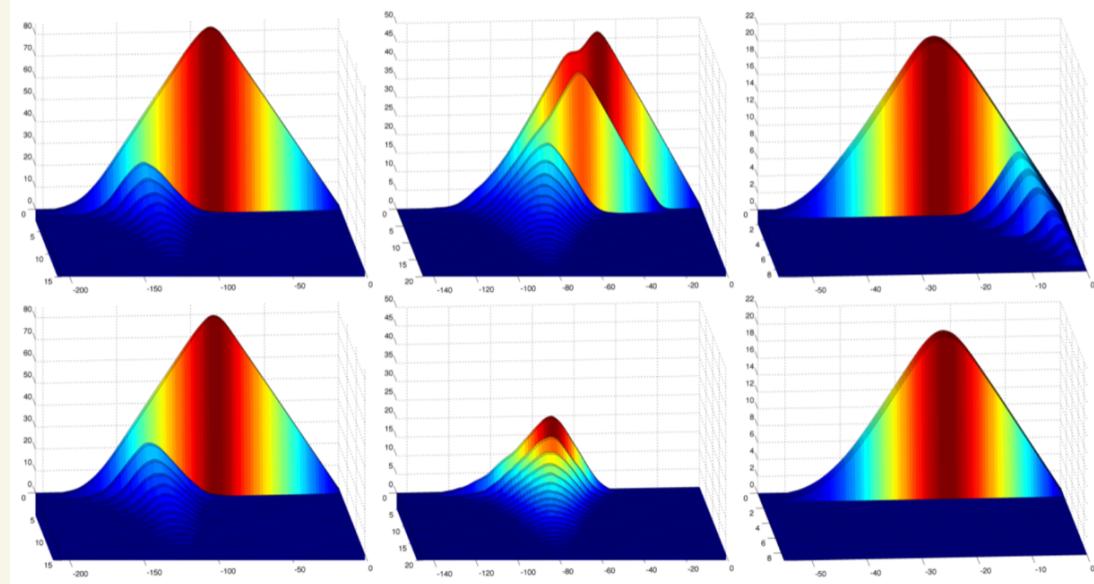
Sphere
Sampled

Build
filtration
(torus here)



Torus →

Sphere →



dim 0 dim 1 dim 2

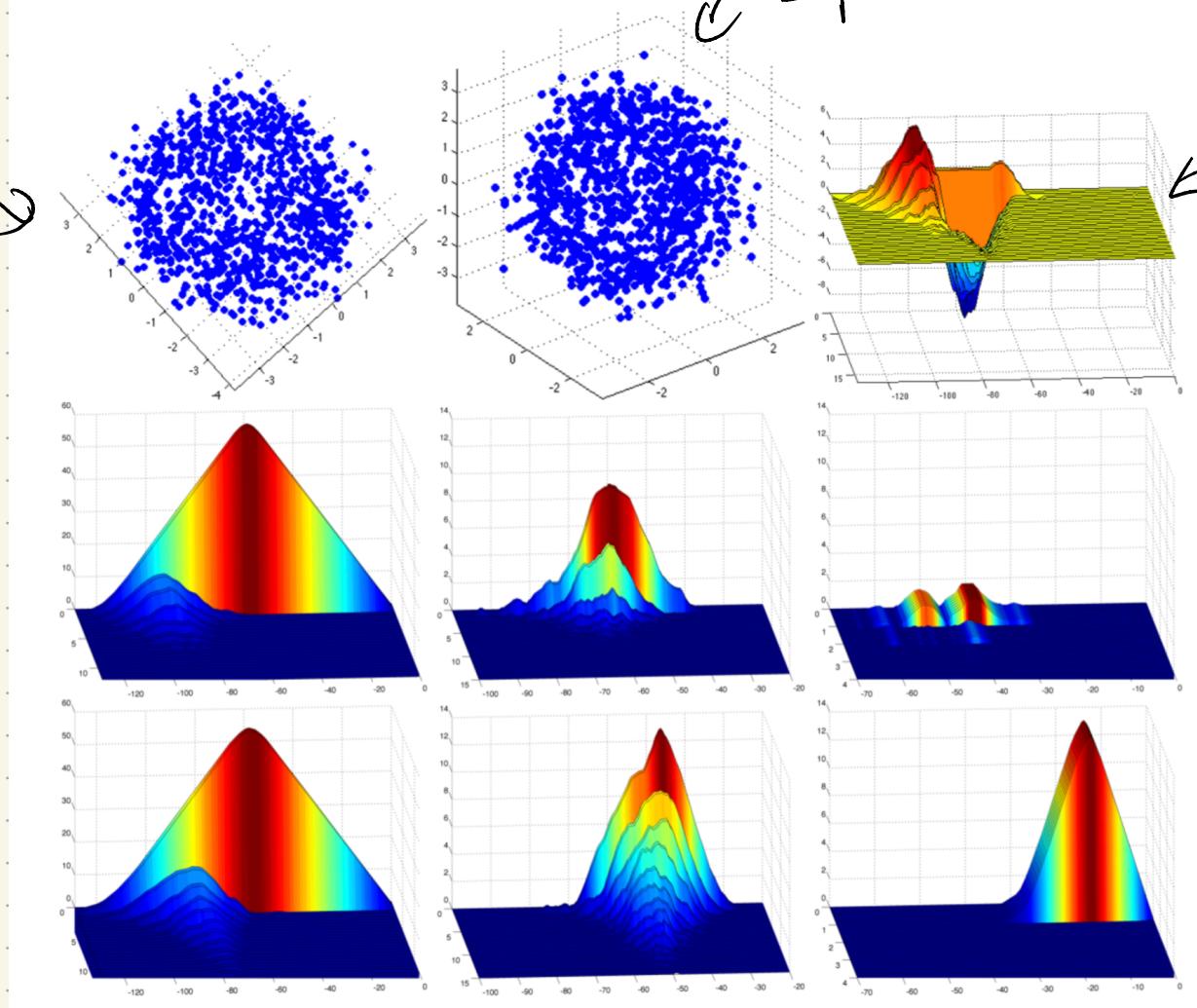
And with noise, can still distinguish!

Sphere

Torus

Torus →

Sphere →



Dim 0

Dim 1

Dim 2

Gaussian
noise

10,000
repetition
average

Pros / Cons

- Implemented in many libraries,
& seems to do well!
 - ↳ can take averages, has notion of stability, & useful in ML
- Average is not itself a Persistence diagram
 - ↳ interpretability then suffers

Barnes, Polanco, Perea 2021

Persistence Images

Adams et al 2017 JMLR

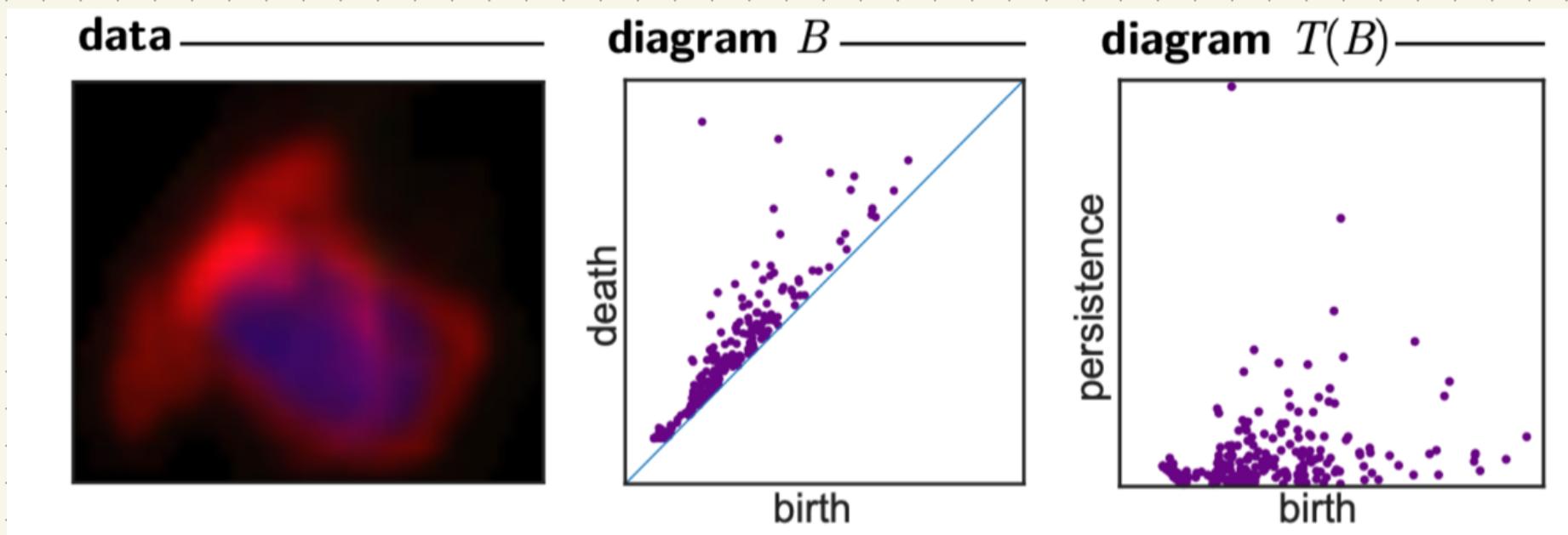
Goal: Represent a persistent diagram
so that:

- Output of representation is a vector $\in \mathbb{R}^n$
- representation is stable & efficient to compute
- interpretable connection to original P.D.
- Can adjust relative importance of points in different parts of P.D.

Trade offs with persistence landscapes:

- PLs are invertible & good averages
- PIs are better for ML

Defining (again via picture) PIs



First, transform diagram B to $T(B)$:

$$T(B) = \{ (x, y-x) \mid (x, y) \in X \}$$

So: $(3, 7) \mapsto$

Now have $(\text{birth}, \text{lifespan})$ pairs

Next, weight the points: for example

$$f(\vec{p}) = \Pr / (\text{max persistence of any point})$$

(or could choose other weighting!)

Add Gaussian around points

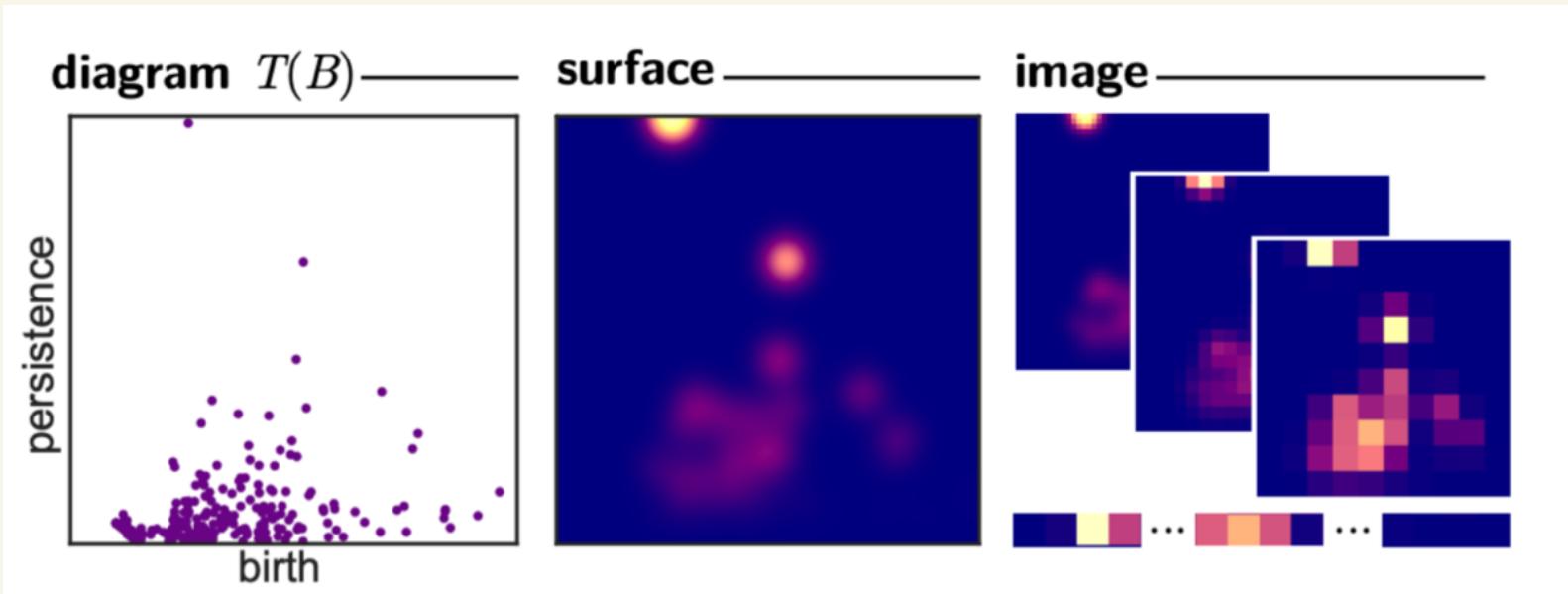
(or some differentiable probability distribution)

$$e_p(z) = \frac{1}{2\pi\sigma^2} e^{-((x-p_x)^2 + (y-p_y)^2)/2\sigma^2}$$

(where σ is usual variance)

Then transform the diagram to
a scalar function on \mathbb{R}^2





Let μ_B

$$\circ \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$x \mapsto \sum_{p \in T(B)} f(p) \cdot e_p(x)$$

This gives a "surface".

Then, reduce to finite dimensional vector
by "boxing" with value = integral of μ
 n boxes

Why?

- Stable!
- Lots of choice in parameters & weights & distribution
- Vector output can just plug & play into most ML pipelines
- Seems useful even with noise, & can infer what region influenced results (to some extent)

Aside:

Code tutorial

→ using GUDHI

