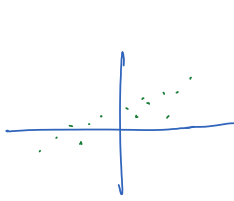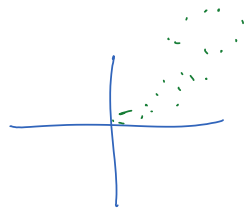# Topological Data Analysis

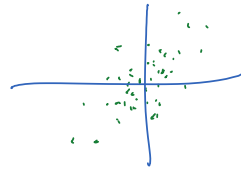Thursday, January 11, 2018    4:13 PM

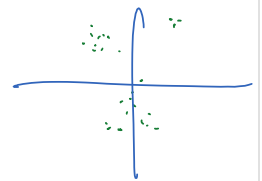If you've done a lot of Data analysis, you know about

linear-ish data       exponential or power-law-ish data       normal (Gaussian) data       clustered data
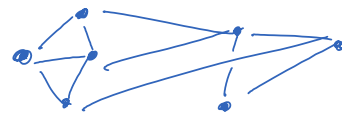
But what about

spatial data

vs

network data

data with a shape

?

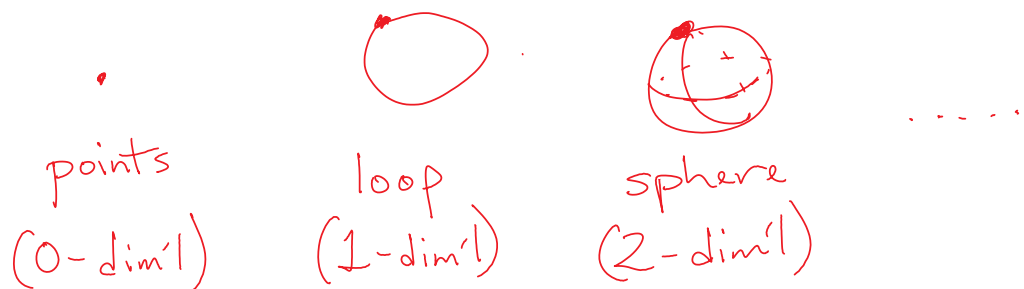Topological data analysis (TDA) allows a fun[damentally] different way to look at this data.
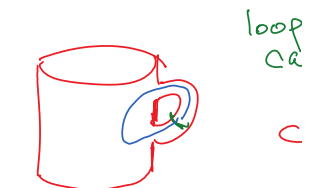
Used fruitfully on
- Cancer data
- diabetes

- ADHD
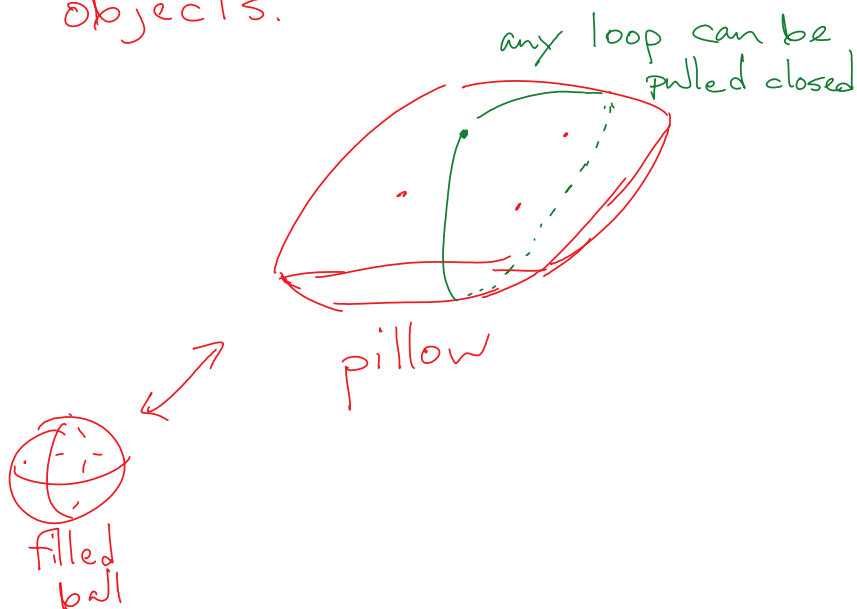- tortuosity of blood vessels
- LiDAR data (ambulance vs Toyota machine gun mount)

# Topology crash course

Pure math version: built on loops, spheres, etc.

points
(0-dim'l)

loop
(1-dim'l)

sphere
(2-dim'l)

. . . . .

These help mathematicians differentiate between objects.

any loop can be pulled closed

pillow

filled ball

loop ca

c

loop

d

blue loops pull closed on the s

The pillow has no way you

can make a loop on the surface
that can't be pulled closed.

Then coffee mug a
each have two a
kinds of loops !
be pulled close
surface.

This pure math "homology"
does not care about distances.
We need a metric (measure, length, distance) for data

[homo ≈ same, so stu
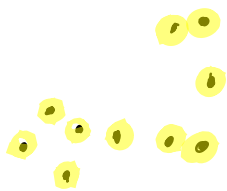objects are "the s
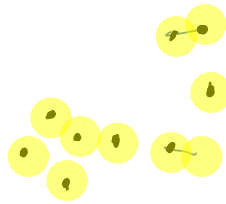
point
(0-dim'l)

line
segment
(1 dim'l)

triangle
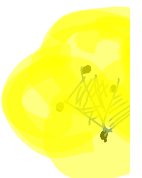(2-dim'l)

tetrahedron
(3-dim'l)

## Vietoris-Rips complex

10 disconnected
points

4 line segments
have appeared

Lots of lines,
and a
few triangles.

Tetr
appe

Grow the radius of a ball around each point. Wh
connected?

Keeping track of how data points get connected to o
data points as you grow balls around them gives y
topological "signature" for your data.
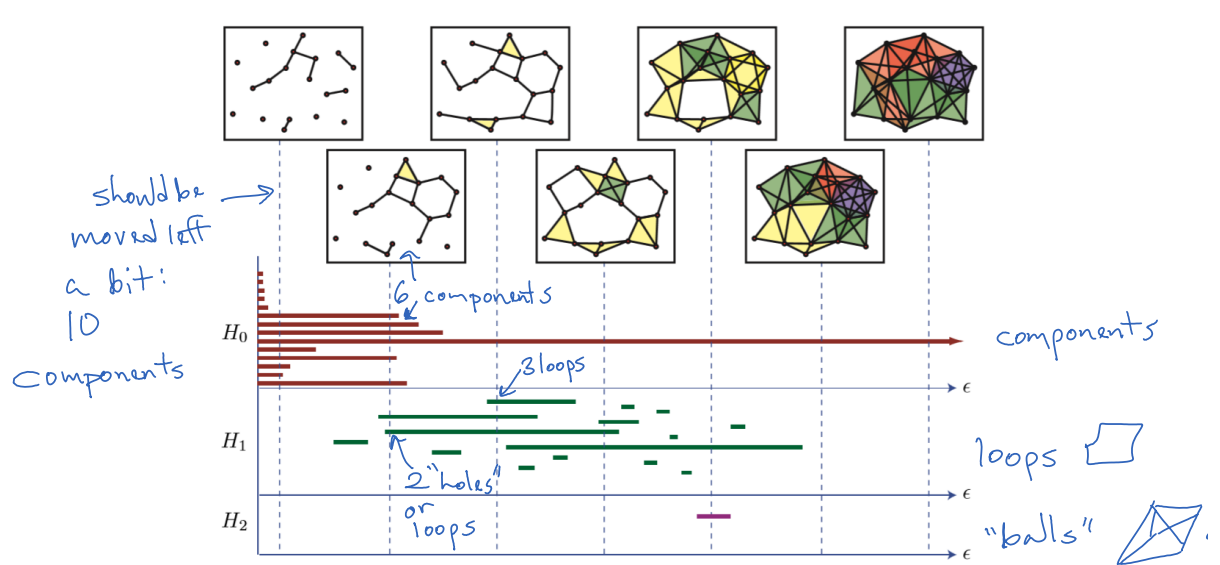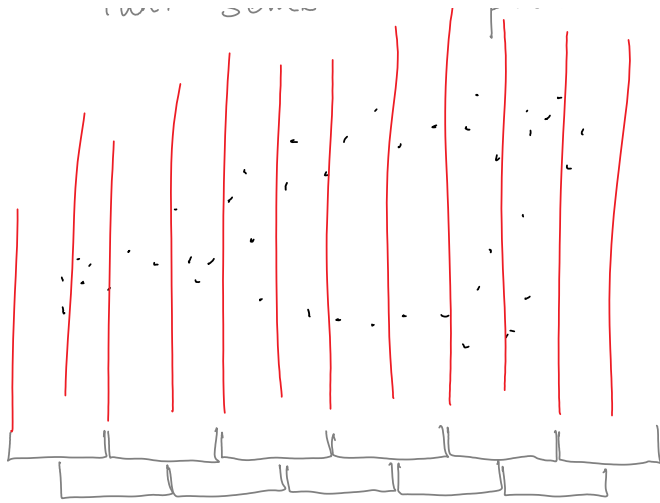
Persistent homology ~~barcode~~ or ~~birth-death~~ dia



should be →
moved left
a bit:
10
components

6 components

$H_0$

components

3 loops

$H_1$

loops

2 "holes"
or
loops

$H_2$

"balls"

FIGURE 4. [bottom] An example of the barcodes for $H_*(\mathrm{R})$ in the
example of Figure 3. [top] The rank of $H_k(\mathcal{R}_{\epsilon_i})$ equals the number
of intervals in the barcode for $H_k(\mathrm{R})$ intersecting the (dashed) line
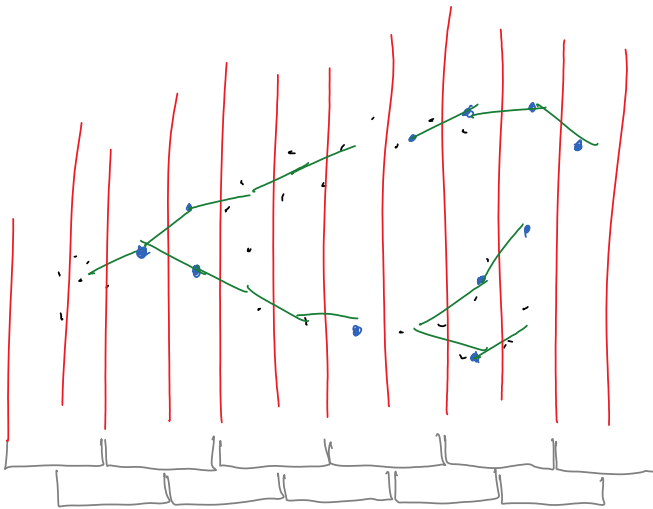$\epsilon = \epsilon_i$.

Screenshot from Rob Ghrist's "Barcodes:
The persistent topology of data."

## Workflow

- Prepare data if necessary (scaling, for instance)
- Choose a way to measure distance
- ( • Project to 1 or 2 dimensions if you want to visua
- Choose a way to slice your data
- Tune some other parameters

slice data like ||||| (as opposed
to ☰ or ◎ or 〝〟… )
with 50% overlap of bins



Once data is slic
clump in each s
and draw edges
there is overlap be
clumps in overlappi

- Check the barcode to s
  if topological features
  persist or if you get
  certain kind of signatu

- Check the graph you g
  against domain experti
  to see if it reveals
  anything!