

Reminders: HW 2 due Friday

Comments on Project proposals:

How to deal with outliers

Density level sets

Metric measure spaces / Wasserstein distance

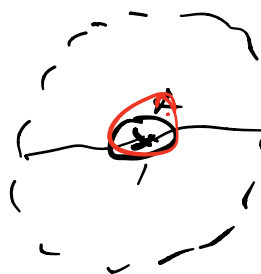
Distance to measure.

$$d_I(P_H(R(X)), P_H(R(Y))) \leq \underline{d_H(X, Y)} \rightarrow \text{A}$$



long H_1
1 H_2

v.s.



H_1 abt half length
short extra H_2

v.s.



long H_1
long extra H_2

Problem: because sampling single point can drastically change P_H .

Common setup: \leftarrow has tails

$$X \stackrel{\text{unif}}{\sim} M + \underline{N(0,1)}$$

$\in \mathbb{R}^d$

\rightarrow some probability of sampling outlier

What we often want is to capture topology of dense subsets, ignore outliers.

E.g. Carlsson & de Silva's idea of using
knn-based co-density filter.

→ what is k ? what is threshold for cut?

one solution for low dimensions:

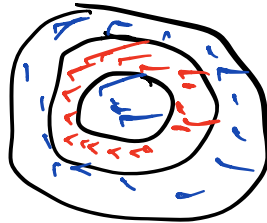
don't do Rips. Do density level sets

$$X \stackrel{\text{unif}}{\sim} S + N(0, \sigma^2)$$

$$F'(a, \infty)$$



pdf:
 \hat{f}



$PH(\hat{f})$: will have
robust fl,
feature

if we have ground truth, this is straightforward.

problems we need to estimate pdf f from
samples. Convergence is exponential in dimension
ie $X \subseteq \mathbb{R}^d$ need $O(\frac{1}{\epsilon}^d)$ samples for
 $O(\epsilon)$ accurate estimator of f .

in low dimensions ($d \in \{2, 3\}$) this is tractable.

can adapt stability results from levelset persistence
 $d_H(PH(f), PH(g)) \leq \|f - g\|_\infty$

so if f is pdf, g is \hat{f} and estimator converges
uniformly as $n \rightarrow \infty$, then barcodes converge.

↗ assume Euclidean.

In higher dimensions ($d \geq 4$) need something else.

Metric measure spaces are good definition to work with.

Def: A metric measure space is a triple

(X, d, μ) where (X, d) is metric space
 X is set, d is metric, μ is measure

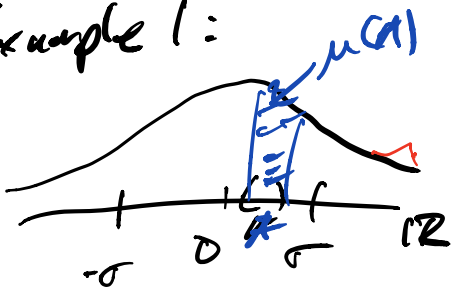
μ is a measure in sense of prob. measure.

mass of Borel subset $A \subseteq X$ is $\mu(A)$

and $A \cap B = \emptyset \Rightarrow \mu(A \cup B) = \mu(A) + \mu(B)$

(countably) additive

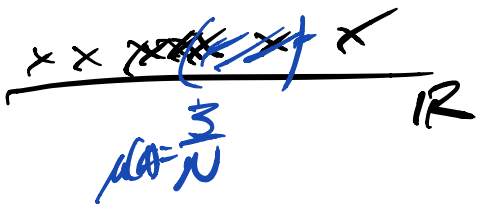
Example 1:



pdf $N(0, \sigma^2) \sim \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2}$

$$\mu_f(A) = \int_A f(x) dx$$

Example 2: $X \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, N samples



empirical measure:

$$\mu_x(A) = \frac{1}{N} |X \cap A|$$

$$= \frac{1}{N} \sum_{x \in A} \{x \in X\}$$

$$= \frac{1}{N} \int_A \sum_x \delta_x(y) dy$$

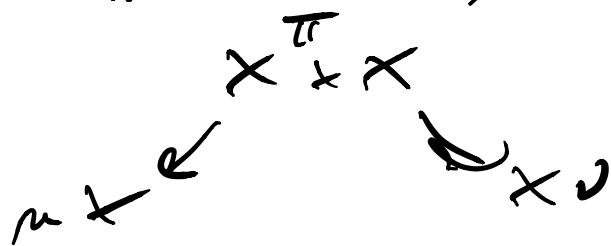
in what sense are measures close? δ_x Dirac delta

i.e. we might expect $\mu_x \rightarrow \mu_f$ if $x \sim f$
 $N \rightarrow \infty$.

Note: positive measures have $\mu(A) \geq 0 \forall A$
 probability measures have $\mu(X) = 1$
 $\mu(\mathbb{R}^k)$

Transport plan: let μ, ν be measures on X ,
 $\mu(X) = \nu(X)$ (same mass). A transport plan
 btw. μ and ν is a measure π on $X \times X$

s.t. $\pi(A \times X) = \mu(A)$, $\pi(X \times B) = \nu(B)$



idea is that all mass
 in μ is transported
 to mass in ν

Example: $\mu(\cdot) = \frac{1}{2}$

in discrete setting, can encode in a

Assignment matrix is. $C = \begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \frac{1}{2} \frac{1}{2} \delta_x$

sum rows: projection onto μ

sum cols: projection onto ν

There are many possible transport plans e.g.

$\begin{pmatrix} \frac{1}{3} & \frac{1}{6} & 0 \\ 0 & \frac{1}{6} & \frac{1}{3} \end{pmatrix}$ also valid $\frac{1}{2}$

k transport plans: $\sum_{i=1}^k w_i C_i$ $\sum_{i=1}^k w_i = 1, w_i \geq 0$

is also a transport plan.

$$c_i \geq 0$$

The p th order cost of a transport plan π

$$\begin{aligned} \text{is } C_p(\pi) &= \left(\int_{x \times x} d(x, y)^p d\pi(x, y) \right)^{1/p} \\ &\leq d(x, y) d\pi(x, y) \end{aligned}$$

Note: $C_1(\pi)$ is "Earth mover's cost": i.e. sum of cost to move 1 unit of mass 1-unit of distance

Let $\Pi(\mu, \nu)$ denote set of all transport plans btw μ and ν .

Def: the p th order Wasserstein distance btw μ and ν is

$$W_p(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} C_p(\pi)$$

If measures have finite p th moments, $W_p(\mu, \nu) < \infty$

Comment 1: There is a generalization of Wasserstein distance to barcodes. W_0 is bottleneck distance.

Comment 2: There is a notion of Gromov-Wasserstein distance,

Comment 3: Finding π which minimizes C_p is subject of optimal transport. Important in applications. e.g. supply logistics.

Distance to measure (DTM)

Let μ be a prob. measure on (X, d)

and $m \in (0, 1]$ be mass parameter.

define distance $d_{\mu, m}$ to measure μ as

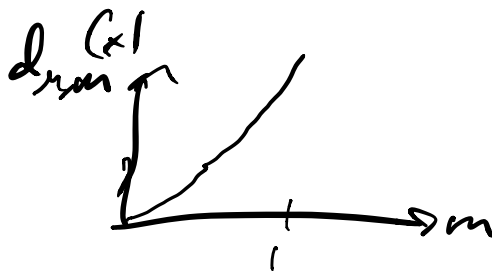
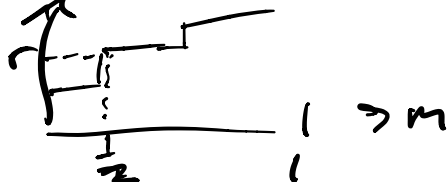
$$d_{\mu, m} : x \in X \mapsto \sqrt{\frac{1}{m} \int_0^m \delta_{\mu, \varepsilon}(x)^2 d\varepsilon}$$

where $\delta_{\mu, m} : x \mapsto \inf_{r \geq 0} \{ \mu(\bar{B}(x; r)) > m \}$



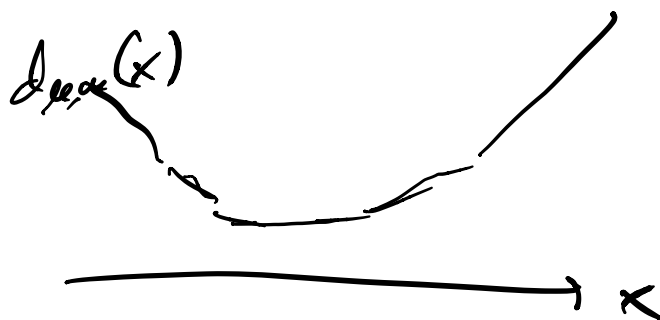
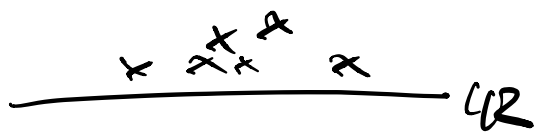
$$\delta_{\mu, \frac{1}{2} - \varepsilon}(x) = \varepsilon_d$$

$$\delta_{\mu, m}(x)$$



for m

μ : empirical measure



↳ points w/ small $d_{\mu, m}$ are "more important" than points w/ large DTM.

• m is a parameter that can be tuned
 $m \rightarrow 0$ distance to closest point
 $m \rightarrow 1$ smooth, may lose structure

idea: create sub-levelset filtration w/
 $d_{\mu, m}$.

Buchet et al 2016 (Oudot Ch. 5.6)
 Thm 3.1:

$$d_{\mathbb{S}}(\text{PH}(d_{\mu, m}), \text{PH}(d_{\nu, m})) = \frac{1}{\sqrt{m}} W_2(\mu, \nu)$$



Proof: ---

A measure ν is a submeasure of μ if
 $\nu(B) \leq \mu(B) \forall B$. Let $\text{Sub}_m(\mu)$ be
 submeasures w/ mass m .

Prop: Let μ be a prob. measure on X , and
 $m \in (0, 1]$. Then

$$\underline{d_{\mu, m}(x)} = \min_{\nu \in \text{Sub}_m(\mu)} \frac{1}{\sqrt{m}} W_2(\mu \delta_x, \nu)$$

↳ mass m
 conc. at x .

Pf: ν is measure of mass $m \Rightarrow \exists$ transport
 plan betw. ν and $m \delta_x$

$$W_2(\mu \delta_x, \nu)^2 = \int_X d_x(y, x)^2 d\nu(y)$$

$$\text{let } d_x = X \rightarrow \mathbb{R} = d_x(\gamma) = d(x, y)$$

$$\text{let } \nu_x : \nu_x(I) = \nu(d_x^{-1}(I)) \quad \forall I \subseteq \mathbb{R}$$

$$\text{let } F_{\nu}(\gamma) = \nu([0, \gamma]) : F_{\nu}^{-1} : m \mapsto \inf \{ t \in \mathbb{R} \mid F_{\nu}(t) > m \}$$

$$\Rightarrow F_{\nu_x}^{-1}(m) = S_{\nu, m}(x)$$

$$\inf \{ t \in \mathbb{R} \mid \nu_x([0, t]) > m \}$$

$$\inf \{ t \in \mathbb{R} \mid \nu(\{ \gamma \mid d(x, \gamma) < t \}) > m \}$$

$$\inf \{ r \mid \mu(\bar{B}(x; r)) > m \}$$

$$S_{\nu, m}(x)$$

$$\int_X d_x(\gamma, x)^2 d\nu(\gamma) = \int_{\mathbb{R}^+} t^2 d\nu_x(t) = \int_0^m F_{\nu_x}^{-1}(l)^2 dl$$

$$\nu \text{ submeasure of } \mu \Rightarrow F_{\nu_x}(t) \leq F_{\mu_x}(t) \quad \forall t > 0$$

$$\Rightarrow W_2(\mu \delta_x, \nu) \stackrel{m}{=} \int_0^m F_{\mu_x}^{-1}(l)^2 dl = \int_0^m S_{\nu, l}(x)^2 dl$$

$$\stackrel{m}{=} \int_0^m \mu(d_{\nu, l}(x))^2 dl$$

inequality is tight for set of submeasures

$$R_{\mu,m}(x) \subseteq \text{Sub}_m(\mu)$$

$$\{ \nu \mid \text{supp}(\nu) \subseteq \overline{B}(x, \delta_{\mu,m}(x)) \}$$

$$\nu(\overline{B}(x, \delta_{\mu,m}(x))) = \mu(\overline{B}(x, \delta_{\mu,m}(x)))$$

Existence of such a measure

$$\mu_{x,m} = \mu(\overline{B}(x, \delta_{\mu,m}(x))), \text{ rescale mass}$$

on $\delta \overline{B}(x, \delta_{\mu,m}(x))$ if too much mass

$$\Rightarrow W_2(m\delta_x, \mu_{x,m})^2 = \frac{1}{m} d_{\mu,m}(x)^2$$

$$\frac{1}{\sqrt{m}} W_2(m\delta_x, \mu_{x,m}) = d_{\mu,m}(x) \quad \square$$

Thm: let μ, ν be prob. measures on (X, d)
mc $(0, 1]$ mass param. Then

$$\|d_{\mu,m} - d_{\nu,m}\|_{\infty} = \frac{1}{\sqrt{m}} W_2(\mu, \nu)$$

$$\text{Pf: } \sqrt{m} d_{\mu,m}(x) = W_2(m\delta_x, \mu_{x,m})$$

Let π be OT plan μ to ν

$$\int_{X \times X} d_X(x,y)^2 d\pi(x,y) = W_2(\mu, \nu)^2$$

Consider submeas $\mu_{x,m}$. $\exists \tilde{\pi}$ submeas. of π
that transports $\mu_{x,m}$ to $\tilde{\nu}$ submeas of ν

$$\Rightarrow W_2(\mu_{x,m}, \tilde{\nu}) \leq W_2(\mu, \nu)$$

$$\text{again, } \sqrt{m} d_{\nu,m}(x) \leq W_2(m\delta_x, \tilde{\nu}) \Rightarrow$$

$$\sqrt{m} d_{\nu,m}(x) \leq W_2(m\delta_x, \tilde{\nu}) \leq \Delta_{\text{ineq.}}$$

$$\begin{aligned} & \leq W_2(m\delta_x, \mu_{x,m}) + W_2(\tilde{\nu}, \mu_{x,m}) \\ & \leq \sqrt{m} d_{\mu,m}(x) + W_2(\nu, \mu) \end{aligned}$$

Similarly,

$$\sqrt{m} d_{\mu,m} \leq \sqrt{m} d_{\nu,m}(x) + W_2(\nu, \mu)$$

$$\Rightarrow \|d_{\mu,m} - d_{\nu,m}\|_{\infty} \leq W_2(\mu, \nu)$$

Now, can apply levelset stability result.

$$d_B(\text{PH}(d_{\mu,m}), \text{PH}(d_{\nu,m})) \leq \|d_{\mu,m} - d_{\nu,m}\|_{\infty} \leq W_2(\mu, \nu)$$

□.



outliers don't change measure much.

measure assigned to outlier is $\frac{1}{n}$

$$W_2(\mu, \nu) \sim \frac{1}{n} r$$

Problem: this function is defined on
all of \mathbb{R}^d . Want simplicial construction.
Order C_1 5.6.