

## ⑦ Intrinsic Dimension Estimation :-

→ Trivial ID estimation using PCA

↓  
Gap in eigenvalue spectrum

↓

$$\text{ID} = \sqrt{\dots}$$

But this is an upper bound

→ Clear exception :- circle as 2d data manifold



$$\text{ID} = \sqrt{2} \text{ but ID} = 1$$

→ ID estimation is very important for estimating density. If we say  $\text{ID} = 2$ , but it is  $= 1$ , density estimation will be wrong.

→ The dimension of fractal is somehow the same thing as ID.

→ Fractals become so comp., they are strongly entangled with chaotic dynamical systems.

→ makes deterministic systems appear stochastic

→ In dynamical system, we have a map

S, that

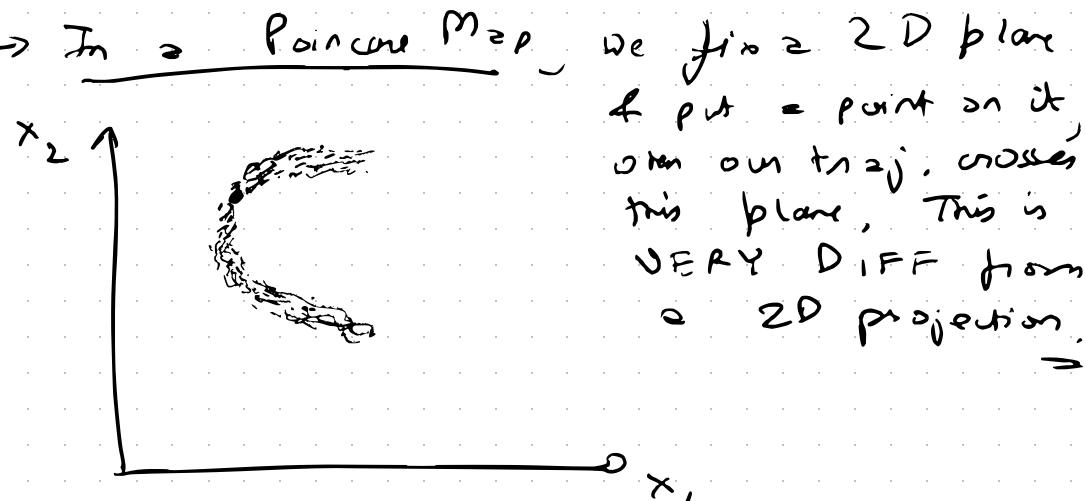
$$x^{t+1} = f(x^t)$$

J

can be deterministic / stochastic

J

$$\rightarrow x^{t+1} = f(x^t | \eta)$$



- Poincaré maps are typically low dim. objects
- our pts. are often confined in a compact region
- first strong hint of the manifold hypothesis.

- Poincaré maps gave us the first hint of the fact that data pts. exists on a manifold that is much lower dim. than the embedding space.
- ID can be non-integer.  
Dimensionality of manifold can also be non-integer / fractal
- So all techniques used to estimate fractal dimension can be used for ID estimation.  
Some of them are :
  - (1) Correlation Dim. Estimator
  - (2) Box - Counting Estimator
- We can estimate ID without any explicit dim. reduc<sup>2</sup>.

## Correlation Dim Estimation :-

→ Count how many data points are at a distance within each other smaller than the threshold  $\rightarrow \{R\}$

$$t_{ijj} \rightarrow \underbrace{\|x^i - x^j\|}_{(\text{Cartesian dist})}$$

How many  $\|x^i - x^j\| \leq R$



$N(R)$

→ monotonically dec func<sup>?</sup> of  $R$

$\therefore$  if  $R \uparrow$ ,  $N(R) \downarrow$

$$\rightarrow N(R) \sim A R^{ID} \rightarrow A \text{ (prefactor)}$$

→ if  $ID = 1$ ,  $N(R)$  grows linearly

$ID = 2$ ,  $N(R)$  grows quadratically

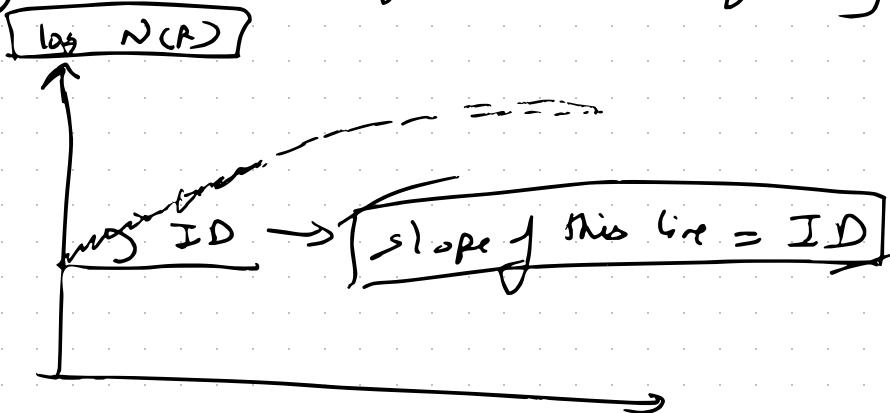
etc.

$$\rightarrow N(R) \sim A R^{ID}$$

↳ defines the correlation dim.

$\rightarrow$  estimate  $N(R)$  from data, & then  
do a  $(\log - \log)$  plot

$\therefore \log N(R) \simeq \log A + ID \log (R)$



$\rightarrow$  Imagine we do the following except -

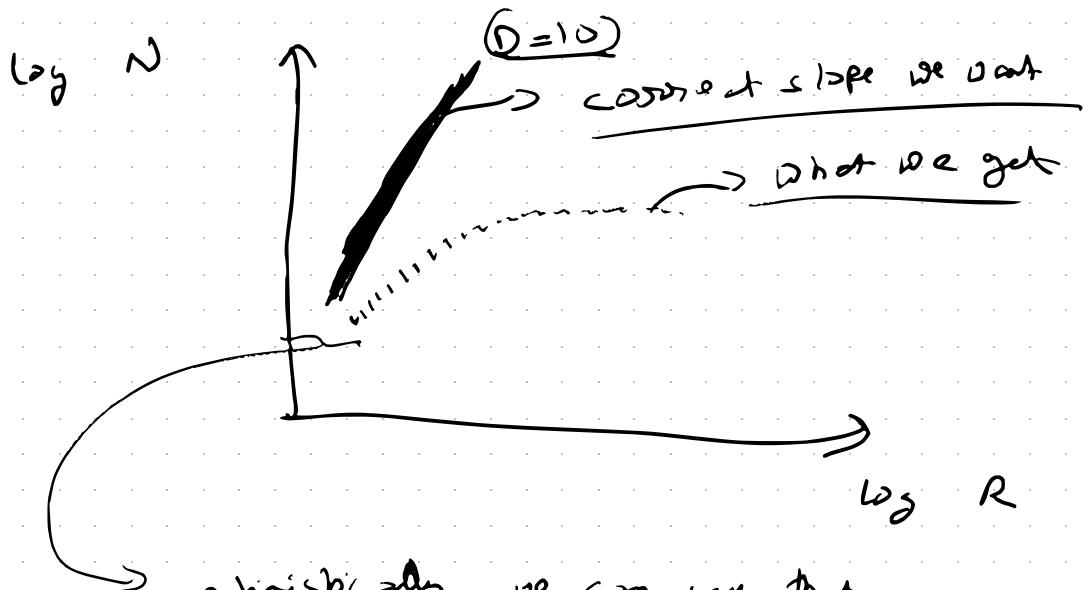
- we harvest data pts in a 10-D space from a Gaussian distrib<sup>2</sup> of variance = 1

$$X \in \mathbb{R}^{10}$$

$$x_e \sim N(0, 1)$$

(Normal distrib<sup>2</sup>)

# of data pts. = 1000 (say)



op timistically, we can say that as

$R \rightarrow \infty$ , slope  $\rightarrow 10$ , but we never really reach it.

$\Rightarrow$  What is the origin of this discrepancy?

prefactor =  $\alpha \rightarrow$  density

$$NCR = \int R^{ID} d\Omega_{ID}$$

(density  $\propto d\Omega \propto R^{2d}$ )

(Vol<sup>m</sup> of unit sphere in ID dimensions)

$$\hookrightarrow \text{in dim=2} \quad \Omega = \pi$$

$$\text{dim=3} \quad \Omega = \frac{4}{3} \pi$$

$$(\pi R^2 \rightarrow R=1 \rightarrow \pi)$$

$$\ln N(R) = \ln \int + ID \ln R + \ln \Omega_{ID}$$

Assumption :- density is const. over the dist.  $R$

$\hookrightarrow$  if density varies, the expression is much more complex. This leads to the discrepancy.

→ It can be shown that the scaling is correct for  $\lim_{\underline{R \rightarrow 0}} R \rightarrow 0$ , i.e. density would be const. in an infinitesimally small neighbourhood.

$$\Rightarrow \boxed{ID = \lim_{R \rightarrow 0} \frac{d \ln(N)}{d \ln R}}$$

(i.e.,  $\log N(R) = \log A + ID \log R$   
is valid only for very small  $R$ )

↳ so in neighbourhood of very small  $R$ , this scaling is exact, even if density is  $\text{pos}^2$  dependent.

- Unfortunately taking  $R \rightarrow 0$  in practice is very hard; ID needs to be very small & we have MANY data points.
- This problem is not very relevant for fractals; they were usually studied in low dim space.
- This led to the origin of other techniques to measure fractal dimension / ID.

→ To recap, origin of problem:

$$\log N \approx \log f + ID \log R$$

$\log N$  scales with density,

so if density varies we're screwed

if we're trying to measure ID &

$f$  & ID are entangled.

⇒ How to disentangle  $f$  &  $ID$ ?

◻ Two - NN estimator :-

→ for every point  $(i)$ , we find  
the first & second nearest  
neighbours.

→ measure the distance,  $r_i$ , of the  
these 2 nearest neighbours to  $\underline{(i)}$ .

$$r_i^{(1)} = \min_{j \neq i} \|x^i - x^j\|$$

$$r_i^{(2)} = \min_{\substack{j \neq i \\ j \neq 1^{\text{st}} \text{ nearest neighbour}}} \|x^i - x^j\|$$

→ For each  $(i)$ , estimate

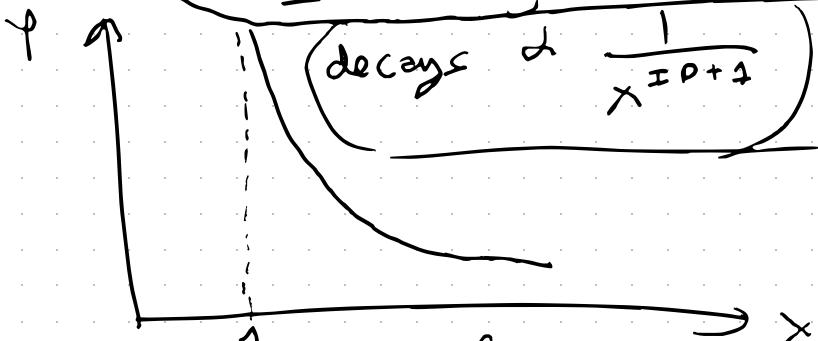
$$\mu_i = \frac{r_i^{(2)}}{r_i^{(1)}} > 1 \quad \text{Calculated by } \underline{\text{deg}^C}$$

→ it can be proved that regardless of ( $f \rightarrow$  density),  $\mu_i$  can be thought of as harvested from a probability density that's a Pareto distrib<sup>2</sup> parameterized by the ID.

$$\mu_i \sim \text{Pareto}(\text{ID})$$

→ Pareto Distrib<sup>2</sup> :- prob. density.

$$f_{\text{ID}}(x) = \begin{cases} \frac{\text{ID.}}{x^{\text{ID}+1}} & x > 1 \\ 0 & x < 1 \end{cases}$$



→ power-law distrib<sup>2</sup> with decay as  $(\text{ID} + 1)$ .

$$\rightarrow \int_{-\infty}^{+\infty} p_{ID}(x) dx = 1$$

non normalized over the density

Pareto distrib<sup>n</sup> should be = 1

→ regardless of  $\mu_i$  can be thought to be generated from Pareto distrib<sup>n</sup>,

$$\mu_i \sim \text{Pareto}(ID) + p$$

(i) Qualitative reason why  $p$  (density) prefactor disappears from the prob.

density :- we consider  $\mu_i$  which is a ratio b/w  $r^{(2)} + r^{(1)}$ .

→ if density is high & data pts. are close by, both  $r^{(1)} + r^{(2)}$  are typically small.

→ if we are in the region of very low density,  $r^{(2)} + r^{(1)}$  are both typically large.

→ ; through dist. scales with density ratio of dist,  $\rightarrow$  such that density prefactor disappears

Ratio of distances ( $\mu_i$ ) scales in a manner which is density independent,

→ ∵ each  $\mu_i \sim \text{Pois}(ID)$   
we can infer ID as described now,

$$\rightarrow P(\mu_i | ID) = \frac{ID}{ID+1} \mu_i$$

(prob. of observing  
 $\mu_i$ , given  $ID$ )

→ for  $N$  data points we have  
several observations  $\{ \mu_i \}$

$$P(\mu_1, \dots, \mu_N | ID)$$

$$= \prod_i \frac{ID}{\mu_i^{ID+1}}$$

assuming each observation  $\mu_i$

$\mu_i$  is indep. Strictly speaking  
this is NOT true, "many

$\mu_i$ 's are of the nearest neighbors  
themselves but for large ( $N$ ) &  
for simplicity's sake we can  
assume this

simple ex. of parameter  
inference given a model,

→ here the only parameters to be  
inferred is ID,

→ Let's use Bayesian Inference.

$$\therefore P(\mu_1, \dots, \mu_N | ID)$$

$$= \prod_i \frac{ID}{\mu_i^{2D+1}}$$

(prior of ID  
≈ const. term)

$$= \frac{P(ID | \bar{\mu}) P(ID)}{}$$

(maximize (P) w.r.t. (ID))

typical procedure for maximum likelihood estimation (MLE)

↓  
Find (ID) which maximizes  $\frac{P(ID | \bar{\mu})}{}$

( $\bar{\mu} = \text{given observations}$ )

→ we do a typical procedure of maximizing  
 $\log P(ID | \bar{\mu})$ , instead of  $P(ID | \bar{\mu})$

$$\therefore L(ID) = \log P(ID | \bar{\mu})$$

$$[\bar{\mu} = \mu_1, \dots, \mu_N]$$

$$= \sum_{i=1}^N \log \frac{ID}{\mu_i}$$

$$= N \log ID - (ID + 1) \sum_{i=1}^N \log \mu_i$$

To maximize this :

$$\frac{\partial L}{\partial ID} = \frac{N}{ID} - \sum_i \log \mu_i = 0$$

$$\therefore ID = \frac{N}{\sum_i \log \mu_i} = \frac{1}{\langle \log \mu \rangle}$$

TL;DR

① compute

$\#_i$

$$1^{st} + 2^{nd} N.N.$$
$$(r_i^{(1)}) \quad (r_i^{(2)})$$

② compute  $\mu_i = \frac{r_i^{(2)}}{r_i^{(1)}} \#_i$

③

$$ID = \frac{1}{\langle \log \mu \rangle}$$

TWO-NN estimator algorithm

① much more numerically robust than correlation-dim estimators, "

prob. distrib<sup>2</sup> of  $\mu$  is indep. of f  
 $f^{-1}(r) \sim (n)^{-1}$  dep. on f,

## TDO - NN algorithm →

- ① Find first & second NN for every data pt. ( $i$ )

$$r_i^{(1)} \quad r_i^{(2)}$$

- ② Compute,  $\mu_i = \frac{r_i^{(2)}}{r_i^{(1)}} \# i$

- ③  $\mu_i \sim \text{Pareto (ID)}$ ; the local density cancels out.

Allows us to disentangle ID & density estimates

$$\text{ID} = \frac{1}{\langle \log \mu \rangle}$$



## $\Rightarrow$ Derivation of 2-NN :-

- (i) Let's say we harvest data points in a region where the density is constant and is equal to ( $\rho$ ).
- (ii) The vol $\cong$  of this region is ( $V$ ).

$P(n | V \rho) \equiv$  Poisson distrib<sup>2</sup>  
 (prob. of harvesting ( $n$ ) points given  $V \rho$ )

$$\therefore P(n | V \rho) = \frac{(\rho V)^n e^{-\rho V}}{n!}$$

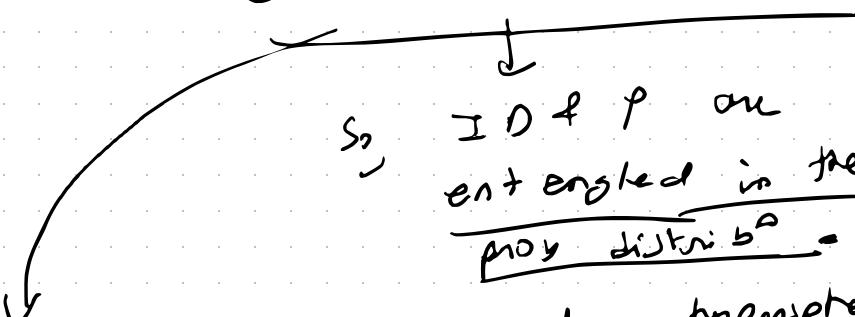
$$\rightarrow n = 0, 1, \dots, \infty$$

we can harvest any # of data pts

$$\sum_{n=0}^{\infty} P(n | V \rho) = 1$$

(Normalization cond?)

- If we fix the size of a region, this tells us how likely it is we observe a given # of pts. in this region.
- All ID / density estimates depend on this prob. distrib<sup>=</sup>.
- The reason  $\rho$  & ID are entangled, is because they appear together in this prob. distrib<sup>=</sup> ( $\therefore \rho = \sqrt{2} ID^{ID}$ )



So, ID &  $\rho$  are entangled in the prob. distrib<sup>=</sup>

$\vartheta$  is always some form of a hypersphere in our derivations.

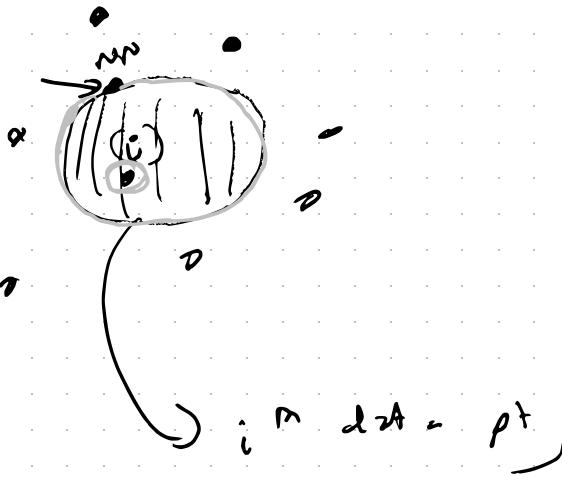
(C) Poisson distrib<sup>=</sup> → prob. of having pts. from a region given the density.

→ Postulate 1 our derivation is based on the fact that fall within a region of vol<sup>c</sup> ( $V$ ) with density ( $f$ ) are harvested acc. to a Poisson distribution.

→ What is the prob. of observing  $\Omega$  data pts. in a region?

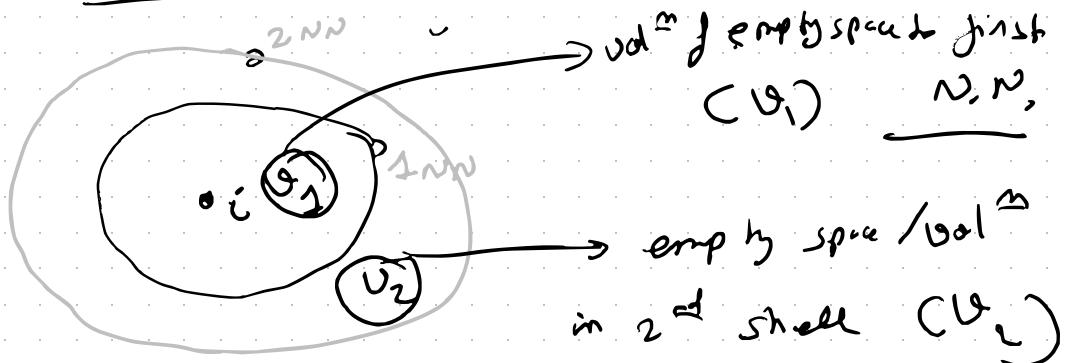
$$P(n=0 \mid \rho V) = e^{-\rho V} = P(\rho V \mid 0)$$

→ Now we can use Bayes' theorem to estimate that density =  $\rho \neq \text{vol}^c = V$  given that we have  $\Omega$  data points in a region. This is the key ingredient of density & ID estimation.  $\Rightarrow$



if draw a sphere with radius dist to NN  
centered at  $i^m dA = p^t$ .

- $\rightarrow$  vol<sup>m</sup> of sphere b/w<sup>2</sup>  $i^m dA = p^t$   
if its NN reduces if density  $T$ .
- $\rightarrow$  Consider 2 diff regions of empty space:



$$P(\ell v = l) = e^{-l} \quad [l = \ell v \\ \text{is written}]$$

$$\rightarrow \int_0^\infty e^{-l} dl = 1$$

(normalization)

$$\rightarrow \text{estimate } P(v|p)$$

prob density  $\int v^a - \text{given } p$   
 $\rightarrow P(x) \rightarrow y = ax \rightarrow \text{we want } P(y)$   
 $\therefore dx \cdot P(x) = dy \cdot P(y)$

$$P(y) = P(x) \frac{dx}{dy}$$

$$P(y) = P\left(\frac{y}{a}\right) \cdot \frac{1}{a}$$

$$\Rightarrow x = y \rightarrow \text{if } y = v \text{ and do the same process we set } P(v|p)$$

$$\therefore P(v|p) = p e^{-pv}$$

$$\rightarrow v_1 \sim \exp(p)$$

$$v_2 \sim \exp(p)$$

$$\rightarrow \text{we want } P\left(\lambda = \frac{v_2}{v_1}\right)$$

$$\therefore P\left(\lambda = \frac{v_2}{v_1}\right) = \int_0^{\infty} dv_1 \int_0^{\infty} dv_2 p e^{-pv_1} p e^{-pv_2}$$

$$\delta\left(\frac{v_2}{v_1} - \lambda\right)$$

How?

delta-func<sup>2</sup>

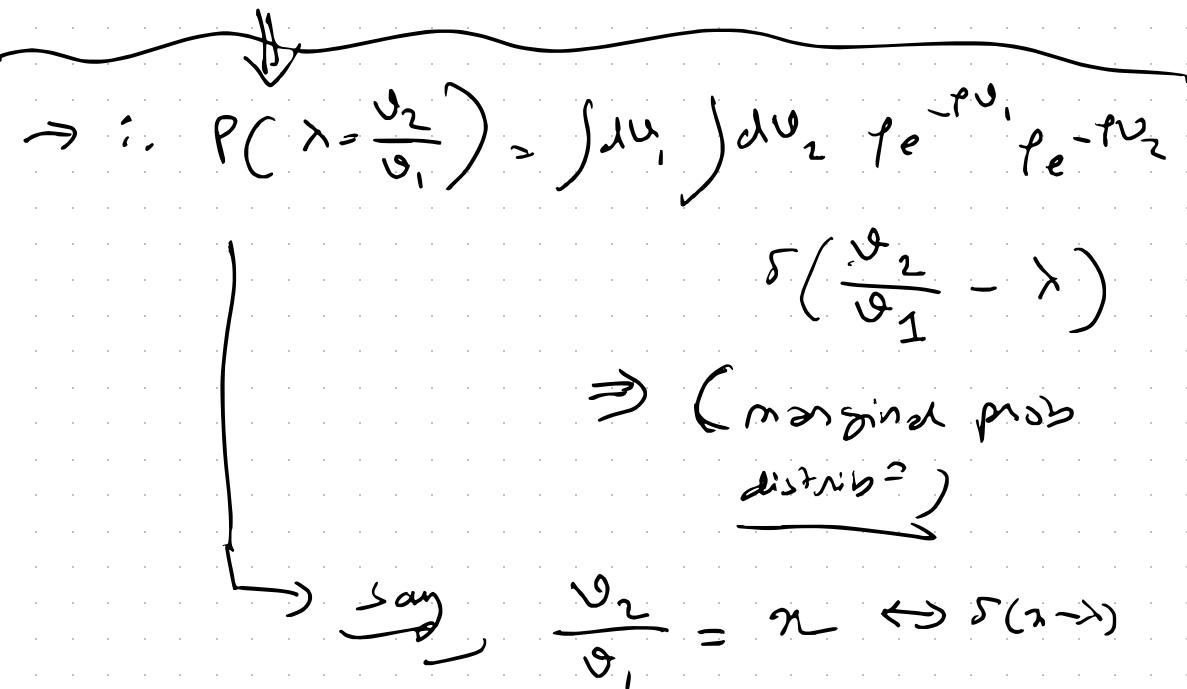
①  $x, y \rightarrow 2 \text{ random vars} \rightarrow \text{generated from } P(x) \text{ prob. density}$

② let's say we have

$$\text{some vars, } s = f(x, y)$$

① how do we estimate  $P(s)$ ?

$$P(s) = \int dx_1 \int dy \quad P(x_1) P(y) \delta(g(x_1, y) - s)$$



$$\therefore P(\lambda) = \int d\lambda \int dx \quad \vartheta_1 e^{-\rho y_1} e^{-\rho x \vartheta_1} \rho^2 \delta(x - \lambda)$$

$$= \frac{1}{1 + \lambda^2}$$

(just integrate over  
 $x_1 \Leftrightarrow \delta$  func<sup>2</sup>  
disappears)

$$\therefore P(X) = \frac{1}{1 + \lambda^2}$$

↓

$\lambda$  has disappeared

while

$P(v_1)$ ,  $P(v_2)$  depend on  $f$

$P\left(\frac{v_2}{v_1}\right)$  is indep. of  $f$ .

This allows us to disentangle ID

& density estimation.

$$\therefore v_1 = \mathbb{E}_{ID} (x^{(1)}) \frac{ID}{C_{vol^3}}$$

hyperplane  
in dim ID]

$$\therefore v_2 \geq S_{ID} \left( (n^{(2)})^{ID} - (n^{(1)})^{ID} \right)$$

$$\therefore \lambda = \frac{v_2}{v_1} = \left( \frac{n^{(2)}}{n^{(1)}} \right)^{ID} - 1$$

$$\therefore \lambda = \mu^{ID} - 1$$

*monotonic func<sup>n</sup> that can be safely inverted*

$$\rightarrow \text{Do know } n \quad P(\lambda) = \frac{1}{1 + \lambda^2}$$

∴ what is  $P(\mu)$ ?

$$P(\mu) = P(\lambda(\mu)) \frac{d\lambda}{d\mu}$$

*; one  
can go  
from  
 $P(\lambda) \rightarrow P(\mu)$*

$$\therefore P(\mu) = \frac{ID}{\mu^{ID} + 1}$$

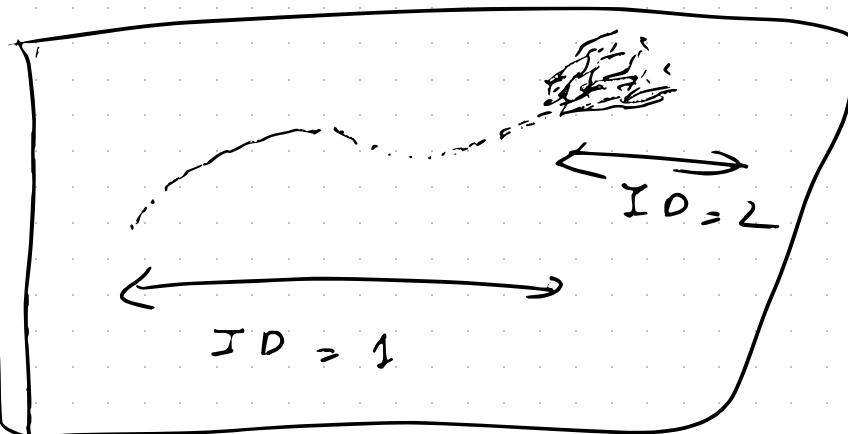
*Parato  
distri b<sup>n</sup>*

∴ This is basically the derivation of  
the TDO-NN estimator.

→ Estimating ID is still a key open  
research problem.

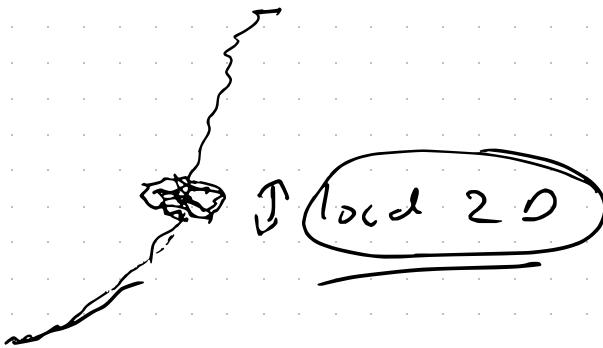
J

∴ ID in diff. regions of our dataset  
can be diff. So pos<sup>c</sup> dep. ID  
estimation is strongly desired.



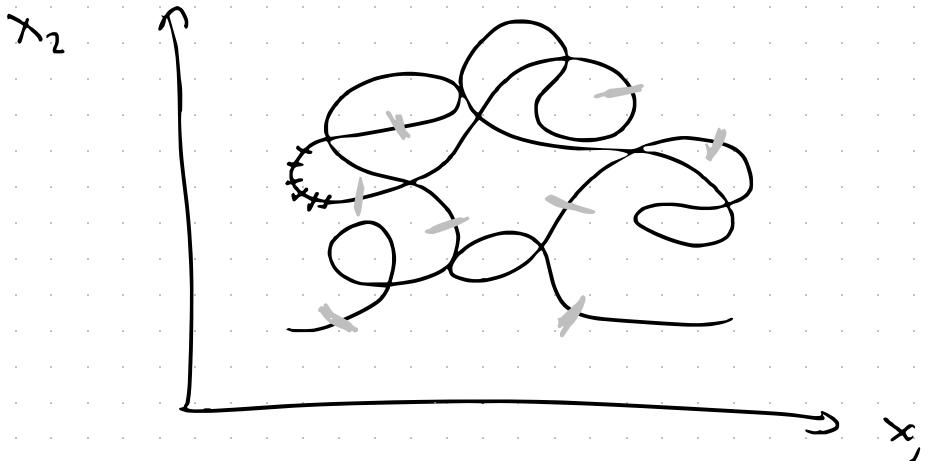
→ Another issue with ID estimation :-

detected pts on the manifold is 2D  
on a really local scale, whereas  
the rest is 1D.



→ There are many more ID estimators  
each with its own advantages &  
disadvantages.

$\Rightarrow$  ID in time-series data:



$\rightarrow$  in ID estimation often we assume Poisson distrib<sup>^</sup>, we assume uncorrelated data. For time-series, data may be correlated.

$\rightarrow$  Let's look at  $(x_2 - x_1)$  phase space traj. from the sim<sup>^</sup> above.

$\rightarrow$  If we honest points very frequently, ID = 1 (if traj. in Hamiltonian

$\left. \begin{array}{l} \text{dynamics, } \& \text{ID} = 1 \times y \\ \text{there's some stochasticity} \end{array} \right\}$   
∴ points are correlated

$\rightarrow$  if we harvest data pts. infrequently  
(grey pts.), in such a way that they  
are approximately independent, then the  
true ID of our manifold emerges

