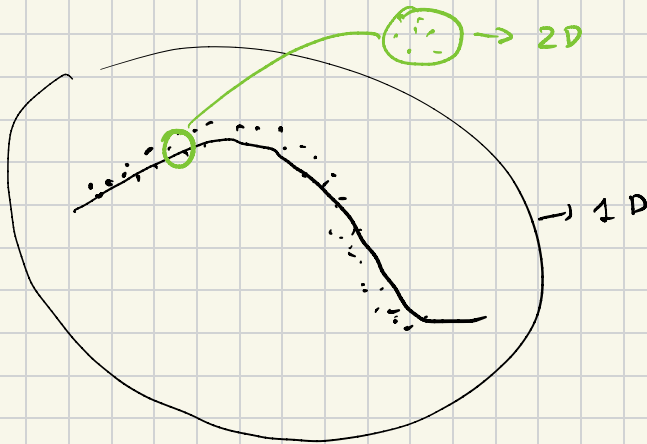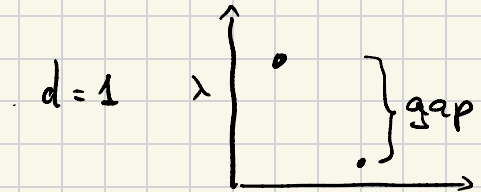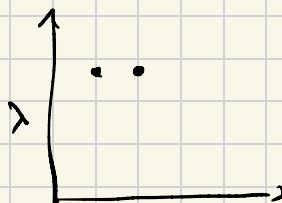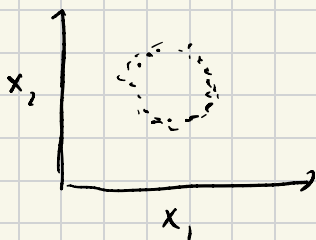# Intrinsic Dim.

a) Minimum Number of variables needed for represent the data with minimum information loss

b) The dimension of the manifold in which our data lies.
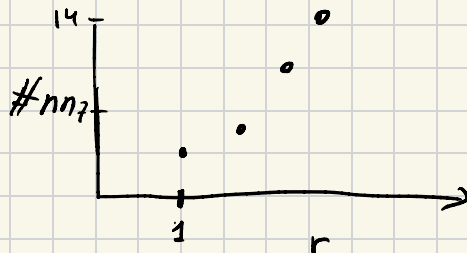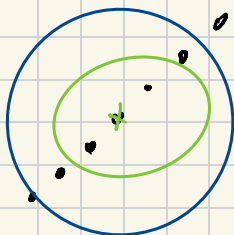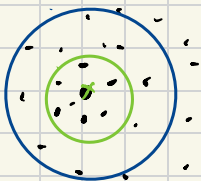


PCA is a method that performs $Id = d$ estimation

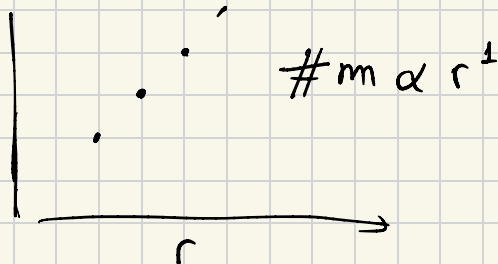$$I_d (PCA) = 2$$

From the data ; learn directly the $I_d$

Fractal



$\#nn$

14

1

$r$

$\#nn \propto r^2$

$\#nn \propto r^1$

$r$

$$\boxed{\# \, nn \, \alpha \, r^d}$$

a) fix a set of different values of r

b) for each data point : count the number of neighbors within r ( take the average)

c) log vs log plot



$log(\#)$ 

$log(r)$

perform a linear fit

d) slope of the fit $\rightarrow Id$
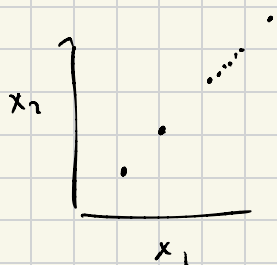
---

For each data point

    a) Compute the distances from its K-NN

    b) Plot the log of the rank as a function of the log of the distances

    c) linear fit

    d) slope will be d

$$log \, (\#) = log \, \rho + d \, log \, (r)$$

count       Id      distances

$$\# = \rho \, r^d$$

We need to disentangle Id & $\rho$ estimations



Two-NN

$$\mu^i = \frac{r_2^i}{r_1^i} \qquad \mu^i = [1, \infty)$$
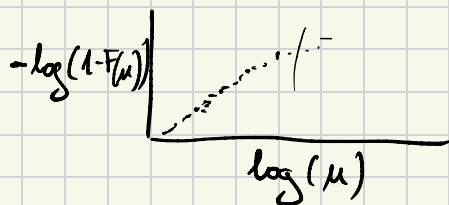
$$\text{i.i.d} \quad p(\mu) = \underbrace{\mu^{-d-1} \cdot d}_{\text{pareto distribution}}$$

$\rho$ within $r_i^i$ can be considered constant

$$\mu^i = \frac{\rho \, r_2^i}{\rho \, r_1^i}$$

ⓐ $$F(\mu) = \int_1^\mu \mu^{*-d-1} \cdot d \; d\mu^* = 1 - \mu^{-d}$$

$$\frac{-\log\left(1 - F(\mu)\right)}{\log(\mu)} = d$$

## ⓑ ML

$$P(\mu) = \mu^{-d-1} \cdot d$$

$$\log \mathcal{L} = \sum \log \left( \mu_i^{-d-1} \cdot d \right) = N \log(d) - (d+1) \cdot \sum_i \log(\mu_i)$$

$$\frac{\partial \log \mathcal{L}}{\partial d} = \frac{N}{d} - \sum_i \log(\mu_i) = \emptyset$$

$$d = \frac{N}{\sum_i \log(\mu_i)}$$