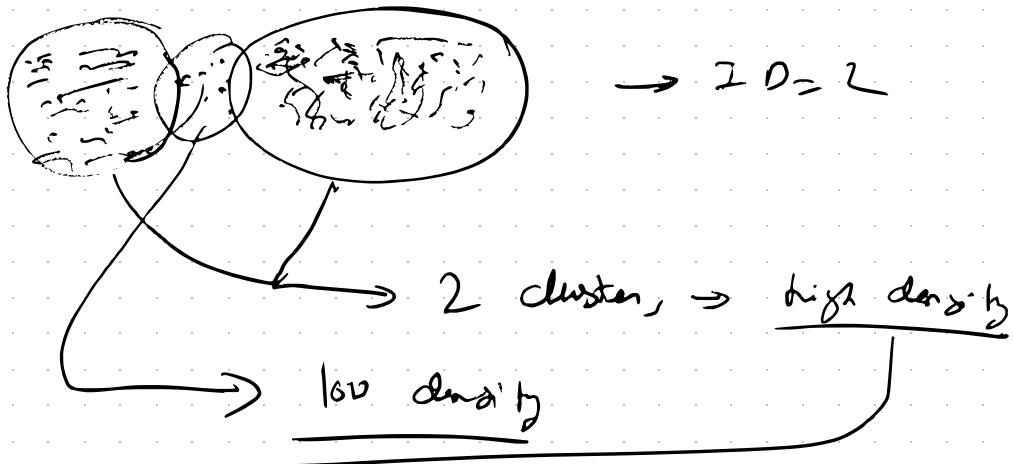


Density Estimation

→ one first estimates ID, & given ID, density is estimated.



→ data typically forms clusters

In molecular systems:

→ high f → metastable states

→ low f → "barrier" regions

regions between = 2 stable states & are explored only rarely. Simulation of rare events.

⇒ in all density estimators we make a compromise in statistical error in estimating (f) & systematic error if (f) is not constant.

"Bias - Variance Tradeoff"

Q. Bias - Variance Tradeoff :-

→ we have performed = preliminary dim. reduce & $\underline{S = ID = 1}$

(let's say, as a simple eg.)

→ So we have only 1 variable.

→ We want to estimate Prob. Density as a func² of 1 variable.

N.B.: Notation related:-

i) $f(x) \equiv$ func² of all coords \equiv Density of data

ii) $P(x) \equiv$ Probability density

① & ② are same except for
normalization.

$$\int dx \varphi(x) = N \quad (\# \text{ of data})$$

$$P(x) = \frac{\varphi(x)}{N}$$

$$\therefore \int dx P(x) = 1$$

→ here we generally take about estimating

$$\varphi(x) \underset{N \rightarrow \infty}{\approx} P(x)$$

→ so density (φ) is a func^c of only 1 variable (y),

→ Data pts. correspond to a set of observations

$$(x_i), i=1 \dots \underline{N}$$

\rightarrow We want to estimate :-

$$\hat{f}^i = f(\hat{Y}^i)$$

f evaluated at \hat{Y}^i

(i) Let's try a naive procedure to estimate

$$\hat{f}^i \underset{i}{\sim}$$

$n_i(\Delta) \rightarrow [\# \text{ data pts. within a dist. } \Delta \text{ from a pt. } (i)]$

data pts. in an interval $\left[Y^i - \frac{\Delta}{2}, Y^i + \frac{\Delta}{2} \right]$

$$\hat{f}^i \underset{\Delta}{\sim}$$

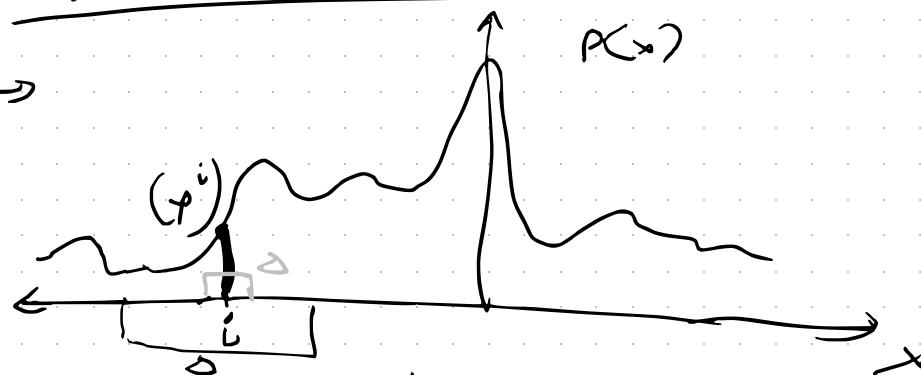


Same procedure as when we make a histogram, "we count # of pts. in each interval"

→ This approach is correct small δ & large N

y δ is small & N is large

① Why & how should δ be small?



① If we want to estimate prob. density at

(i) we want to find p_i , \therefore

δ cannot be large (black range)

& must instead be small (grey range)

→ more rigorously

$$p_i = \frac{1}{\delta}$$

$$\frac{p_i(\delta)}{\delta} \sim \frac{1}{\delta} \int_{y^i - \frac{\delta}{2}}^{y^i + \frac{\delta}{2}} dy \quad p(y)$$

→ in the limit of large N this should be accurate.

$$\frac{n_i(\omega)}{\Delta} \sim \left[\frac{1}{\Delta} \int_{y^i - \frac{\Delta}{2}}^{y^i + \frac{\Delta}{2}} dy \right]$$

$$P(y) \neq P(y^i)$$

$$= \frac{1}{\Delta} \int dy [P(y^i)]$$

$$+ (y - y^i) P'(y^i)$$
$$+ \frac{1}{2} (y - y^i)^2 P''(y^i)$$

∴ we have a correction factor.

∴ we Taylor expand the integral around y^i

cancel out
after
integration

due to parity

$$(\sim \Delta^2)$$

all odd power terms cancel out

$$\therefore \frac{n_i(\Delta)}{\Delta} \sim p^i + \left[\frac{1}{24} p''(y^i) \Delta^2 + O(\Delta^4) \right]$$

\sum : odd contributions
cancel out

estimating the "bias"

a systematic error induced by the fact that
the density is NOT constant.

$\rightarrow p''(y^i)$ is unknown. Only thing
we know is to choose Δ small enough

s. that

$$\frac{1}{24} p'' \Delta^2 \ll p^i$$

very difficult to impose this cond =
rigorously,

→ we take (σ) to be as small as possible, However σ can't be infinitesimally small.

→ if σ is too small, N^i becomes very small (maybe even = 1)

then $\left(\frac{N^i}{\sigma}\right)$ becomes affected by

large statistical error:



$$\frac{N^i}{\sigma} \sim p^i \pm \boxed{\text{error}}$$



if σ is \leftarrow it's very large if
 σ is small

very few observations & if we have too few observations, the error on any estimate becomes very large.

② Estimating the Statistical Errors :

→ "variance" of our estimator.

→ Recall : $n \sim \text{Poisson}(\lambda)$

$$\uparrow \quad \downarrow$$

(n_i)

P

$$\text{in } t \text{ do } \theta = 0$$

of dots ≥ 10 in our interval

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

Some \Rightarrow
values

Step the 2 members using Bayes' formula

↓ (normalization factor)

$$\therefore P(\theta | n) = c \theta^n e^{-\theta}$$

$$\Rightarrow E(\theta) = \frac{n}{d} \rightarrow \text{(expectation value of } \theta)$$

estimates density

$$\text{Var}(\hat{\rho}) = E(\hat{\rho}^2) - (E(\hat{\rho}))^2$$

$$\therefore \text{Var}(\hat{\rho}) = \frac{n}{\Delta^2} \quad \underbrace{(\text{simple algebra})}_{\downarrow}$$

\int
 $\rho(\rho|n)$ func =

$$\therefore \epsilon_{\rho} = \sqrt{\text{Var}(\hat{\rho})} = \frac{\sqrt{n}}{\Delta}$$

(Statistical error on density)

→ "density is not normalized, relative error makes more sense."

$$\therefore \frac{\epsilon_{\rho}}{\rho} = \frac{\sqrt{n}}{\Delta \cdot \rho} = \frac{\sqrt{n}}{\Delta} \cdot \frac{n}{\Delta} = \frac{1}{\Delta}$$

$$\therefore \frac{\epsilon_{\rho}}{\rho} = \frac{1}{\sqrt{n}} \Rightarrow \text{relative error}$$

\rightarrow So for Poisson statistics

relative errors $\propto \frac{1}{\sqrt{\text{no. of observations}}}$



To have small relative errors in the density estimate (f) we must have large (n) .

But, to have large n we must take large (σ) , $\{ \because n \approx f, \sigma \}$

\downarrow
So to reduce Variance, we need large σ

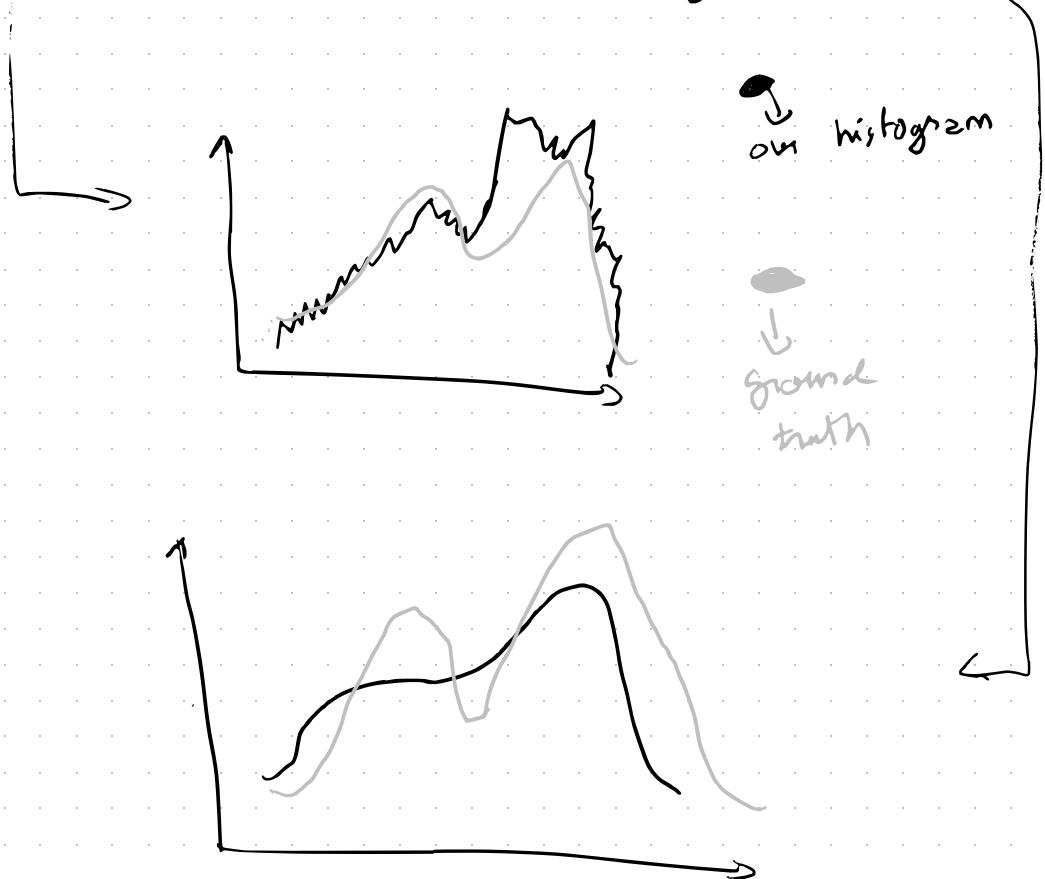
As we saw before, to reduce bias,
we need small σ ,

\Rightarrow So we need to choose σ approximately
to compromise between statistical (variance) &
systematic (bias) errors.

\Rightarrow Some principle as what we do when
making a histogram

\downarrow
We don't want a very small interval,
 \therefore results are very noisy.

\downarrow
We don't want a very large interval
 \therefore then we miss some features



→ This is the origin of the "bias-variance tradeoff"

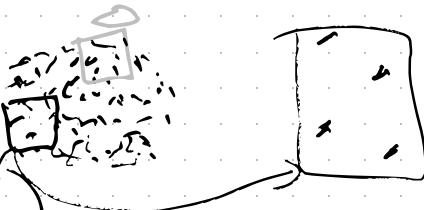
• How to make our density

estimate

adaptive?

What is adaptive?

(too large)



(too small)

→ if we have regions of our data with widely varying densities. Eg molecular systems in beginning of this lecture.

→ If we use a constant Δ for one region Δ is too large, for one it's too small.

\rightarrow So, correct approach is to take
small Δ for high density regions &
large Δ " low " "

$\Rightarrow \therefore$ The smaller the density the larger
the Δ . [the block boxes have
right choice of Δ]

\Rightarrow In the block boxes, the same # of
data pts within Δ is taken,
which means that statistical error is
const

\rightarrow If Δ is uniform, & # of pts.
in it varies widely, some regions
have high statistical errors & others
" low " " ", This is
NOT a good choice, & it is much better
to have const. errors,

→ This philosophy is at the basis of the $K\text{-NN}$ density estimators.

$K\text{-NN}$ Density estimators

1) choose $K = 10$ (say)

2) make σ adaptive. choose
 $\sigma_i = \text{distance of } k^{\frac{m}{m+n}}$

here $\sigma_i = n^{(m/n)^{1/m}} = n^{(10)^{-1/m}}$



$\therefore \sigma_i$ is now point-dependent



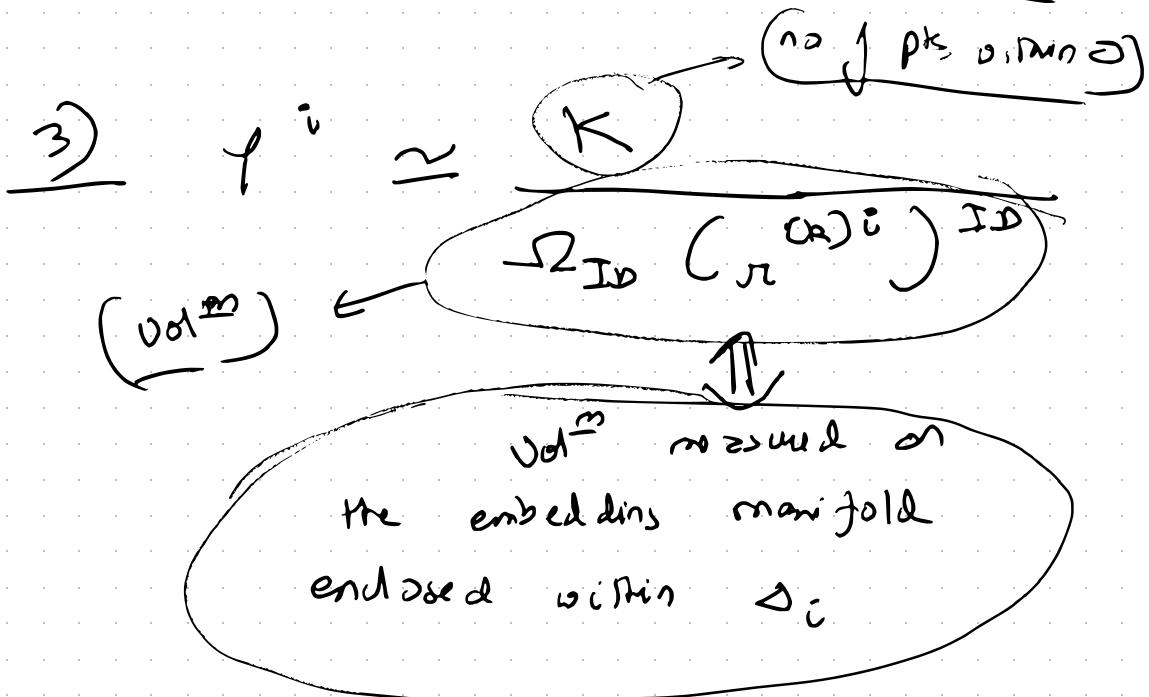
$\therefore K\text{-NN}$ density estimator is adaptive.

(1) Assumption: density is const. on

The scale of σ_i

→ 'y) If we take (k) as small, we can better satisfy this assumption but then σ^2 will have a larger statistical error.

i, here, we have "bias-variance" tradeoff. But the advantage is that "variance" is const, every time.



\Rightarrow 1. This is why estimating ID is vital before computing ID.

if we choose $K = 10$ +

estimate ID = $30 \frac{1}{K}$

we have 10^{30000} which is
ridiculous.

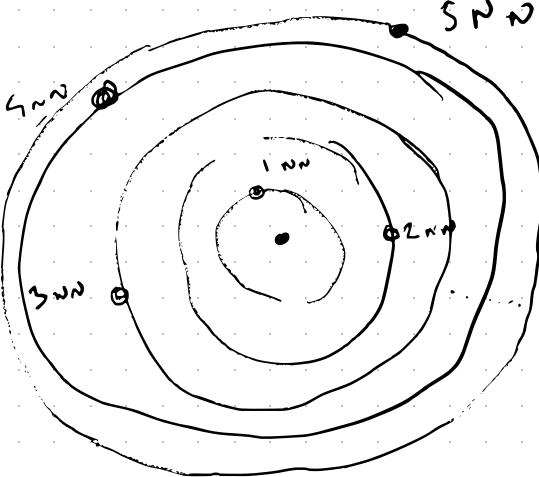
(for numerical stability)

□ Derive the φ_i for kNN estimates

in the form of MLE (Maximum Likelihood Estimation) :-

→ (ID)

→ it should be similar to before,
as the point is that all these
concepts are closely connected.



the various neighbours
+ the $v_i^{(m)}$ from i^m
central pt. to them

\rightarrow say, $k = 5$

\rightarrow each of the $v_i^{(m)}$ of these hyperplanes
 $v_i^{(m)}$ are generated from an exp.
distribⁿ of (φ) density parameter

$$\therefore v_i \sim \exp(\varphi)$$

$$\rightarrow P(v_i | \varphi) = \varphi \exp(-\varphi v_i)$$

Also, $v_i = \sum_{j=0}^D \left[(r^{(j)})^{I^j} - (r^{(j)})^{L^j} \right]$

with $r^{(j)} = O(\log \deg^2)$

∴ What is prob of observing $\{v\}$ given successive vol \equiv observations (v_1, v_2, \dots, v_k)

$$\frac{v_k}{\downarrow}$$

between (t_k) & (t_{k-1}) m no.

$$\begin{aligned}
 & P(\{v\} | v_1, \dots, v_k) \\
 &= \prod_{l=1}^k p e^{-p v_l} \\
 &= p^k e^{-p \sum_{l=1}^k v_l} \\
 &= p^k e^{-p \sum_{l=1}^k (r^{(k)})^{2l}}
 \end{aligned}$$

$\Rightarrow P$ depends on (p) & (r) . So both can't be estimated simultaneously,
However, we have already estimated the (r) .

\therefore We do MLE to compute (\hat{P})

$$\rightarrow \hat{P}: \frac{\partial \log P}{\partial \hat{P}} = 0$$

$$\rightarrow \log P = k \log \hat{P} - \hat{P} \sum_{i=0}^k (n^{(k)})^{ID}$$

$$\rightarrow \frac{\partial \log P}{\partial \hat{P}} = \frac{k}{\hat{P}} - \sum_{i=0}^k (n^{(k)})^{ID} = 0$$

$$\therefore \hat{P} = \frac{k}{\sum_{i=0}^k (n^{(k)})^{ID}}$$

Thus we obtained the eq² we wanted.

N.B. \rightarrow for simplicity's sake, the index

(i) from $n^{(k)i}$ has been dropped to make the notation simpler.

→ Remember :-

- (i) everything is based on the Poisson distribⁿ. In that distribⁿ, the 2 parameters that characterize the data distribution ($\lambda D + \rho$) enter in an entangled manner. But, it can be disentangled by first estimating λD (such as using 2NN) & then given λD we estimate (ρ).



This is particularly relevant for real world data which is high dimensional but where λD is relatively low.

