# Kernal density estimation :- (KDE)

* KDE is a non parametric way to estimate the probability density function of a random variable.

* Non parametric way means the distribution will not follow any specified parameter such as mean, median etc---. simply defined as distribution - free way.

* KDE is a fundamental data smoothing problem where inferences abt the population were made based on a finite data sample.

## what is Kernal ?

* Kernal of PDFs or PMFs is the form of PDFs and PMFs in which any factors that are not functions of any of the Variables in the domain are omitted.

* Simply it means in PDFs the factors which doesn't involved the domain that is the Variable will removed because it will create unnecessary PDFs.

* Kernal is a weighing function

PDFs of normal distribution: $\dfrac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-M)^2}{2\sigma^2}}$

Its associated kernal is
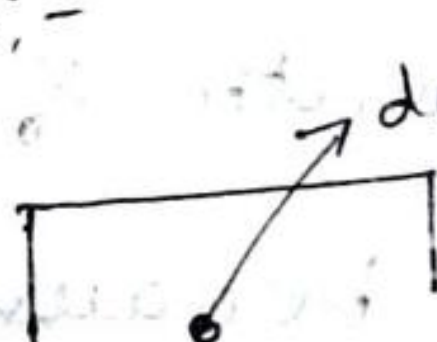
$$P(x|M;\sigma^2) \propto e^{-\frac{(x-M)^2}{2\sigma^2}}$$

Here $\dfrac{1}{\sqrt{2\pi\sigma^2}}$ is removed ~~bee~~ even though it has $\sigma^2$, ~~it ser~~ because it is not the function of $x$.

* A kernal is a non-negative real valued integerable function K.

& ~~It L~~ K is symetric $K(x) = K(-x)$

Some Kernal functions:-

$K(\omega) = 1/2$      Box car

$K(u) = (1-|u|)$      Triangular

$K(u) = \dfrac{3}{4}(1-u^2)$      Epanechnikov

$K(u) = \dfrac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}u^2}$      Gaussian

and many more

## Graphical Explanation with respect to data

* Suppose there are 5 data points in the dataset (just for example).

* on every data point $x_i$, we place a Kernal function $K$
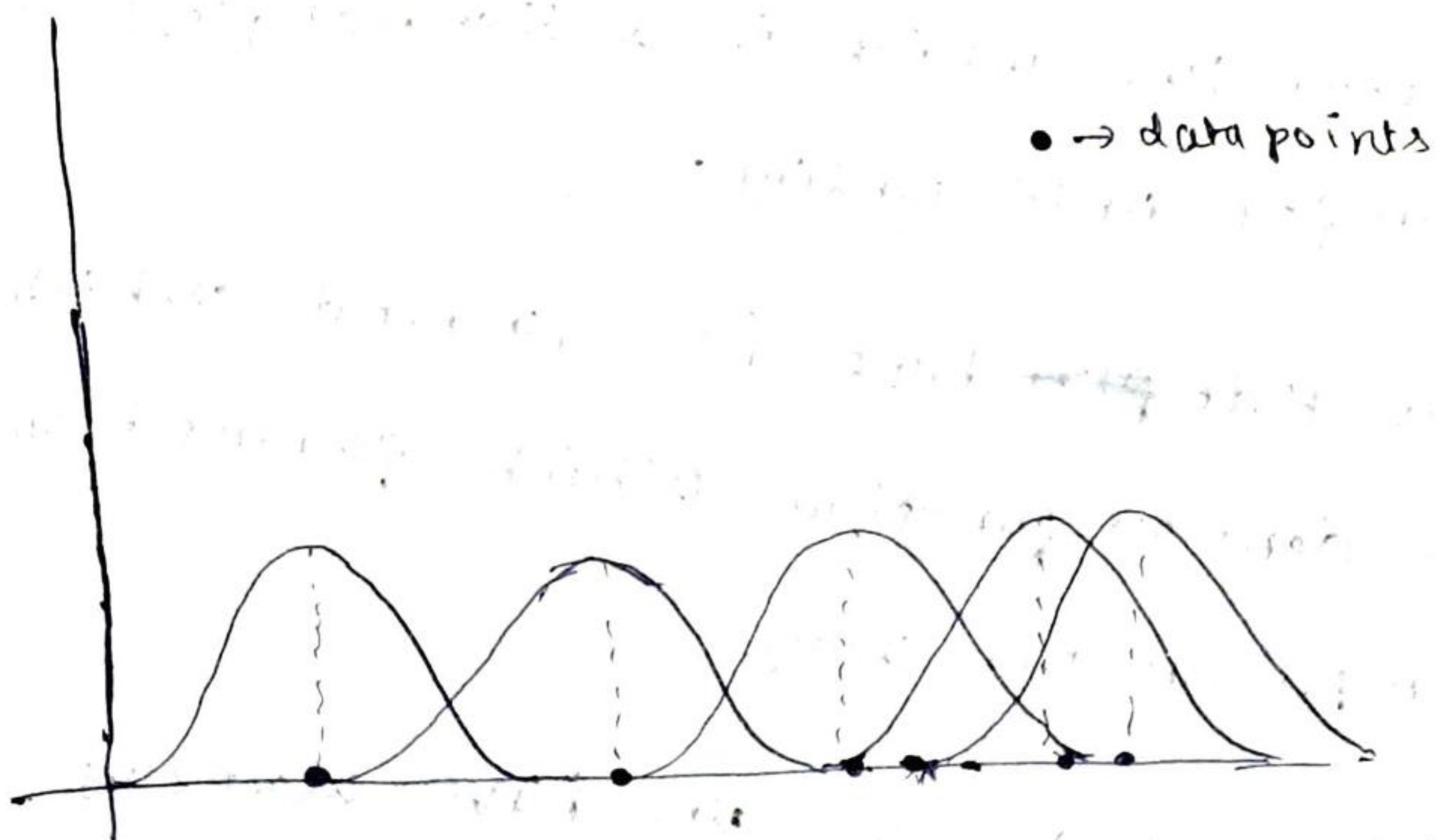
* The Kernal density estimate is

$$\hat{f}(x) = \frac{1}{N} \sum_{i=1}^{N} K(x - x_i)$$

* Since there are more Kernal functions like Boxcar, gaussian, triangular etc..., the choosing of the Kernal is not crucial.

* But choosing of bandwidth of Kernal is more important and crucial because it determines the shape and smoothness of the Kde.
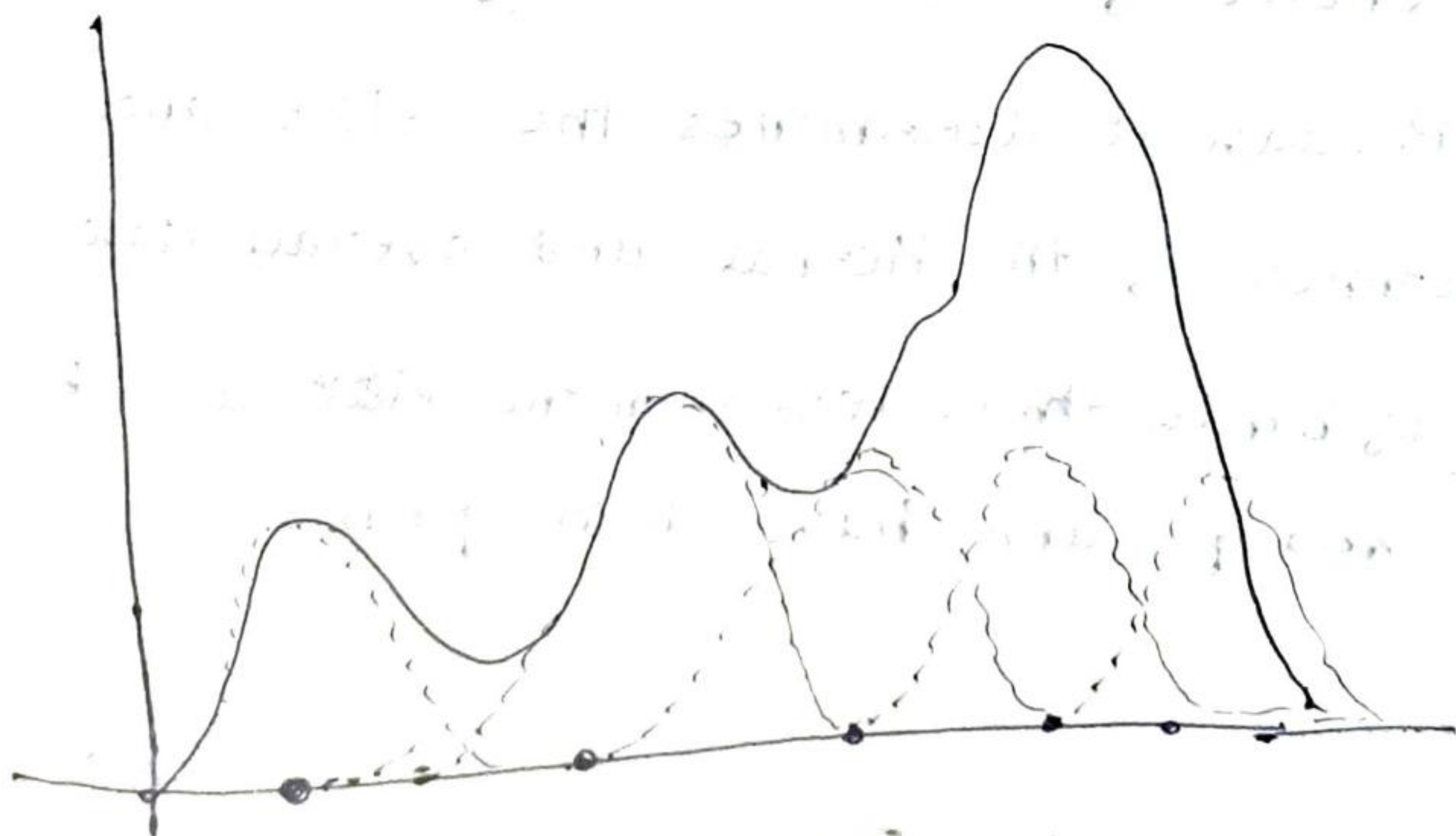
* ~~choos~~ But we don't know exactly for ~~reason~~ what reason choosing of Kernal is not crucial.

• → data points

y Since i have said for example 5 data points I have centered a gaussian Kernal function for each data point. The Kernal have some bandwidth h.

or Kernal functions

y Now it as all the Kernals should be joined are added upto form a Kernal density estimate



y The dark blue line is a kde plot for a feature.

* dotten line which is a Kernal, just shown for understanding.

* The Kde plot Plot is formed which is a density function which formed, in a non parametric way.

* This kind of plot we have already seen in distplot (Histogram + Kde plot).

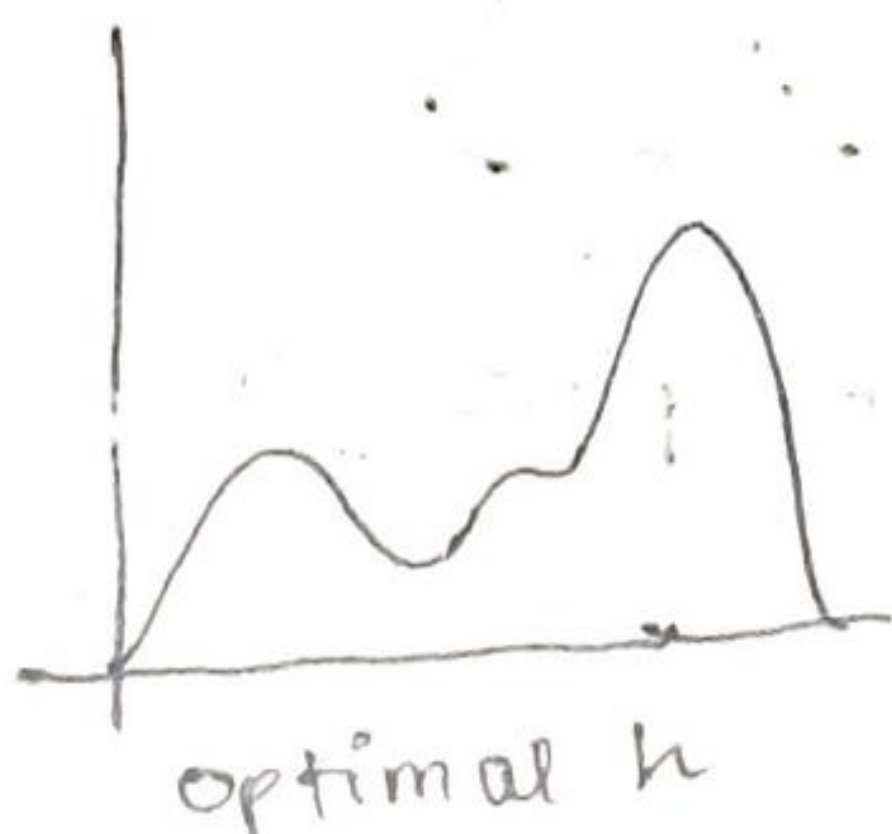* Simple if we join the edges or centres of Histogram we will get Kde plot.

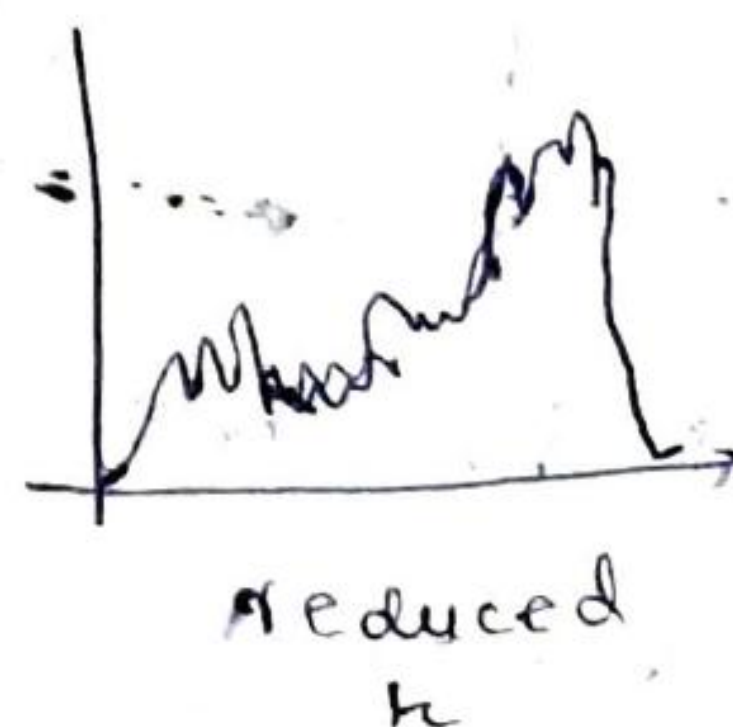* It is mainly Known for its smoothness.

## Choice of bandwidth of Kernals :-

* choice of bandwidth is very crucial point.

* Because it determines the size and smoothness of the Kernal and overall Kde.

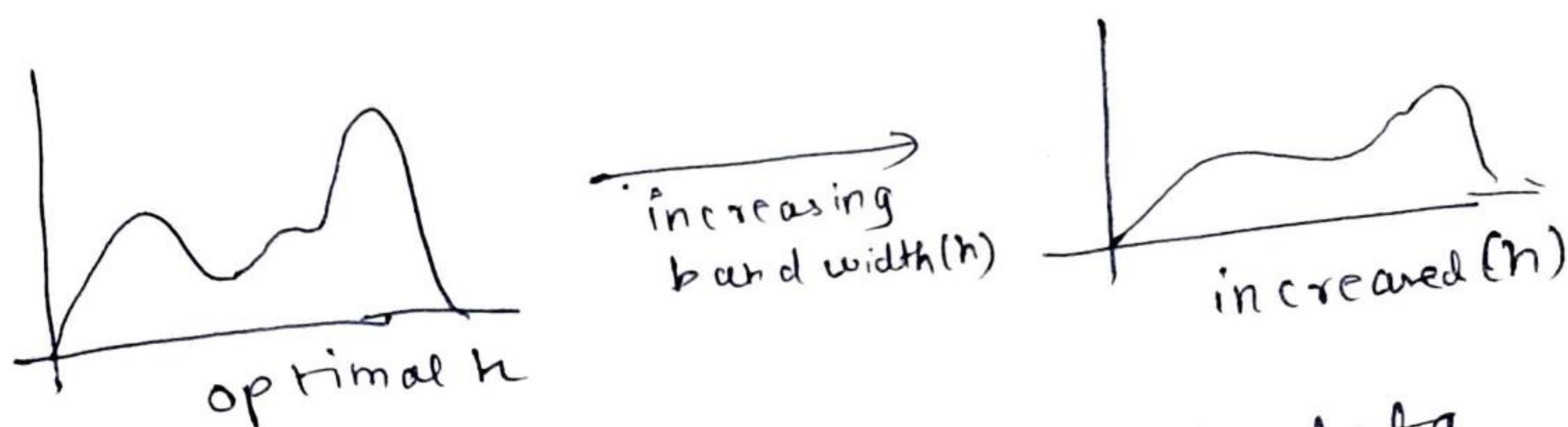* If bandwidth is less then the Kde will be more sharps and have more peaks



reducing
the bandwidth
(h)

optimal h

reduced
h

* If I readuce the bandwidth if ~~they~~ will increase variance of kde plot.

* If I increase the bandwith it will smooth the plot



optimal h

increasing band width(h)

increaved (h)

* If I increase the bandwidth the data may loose its significance because due ~~to~~ to the ~~high but~~ increased bandwidth the data~~t~~ may loose its modality since its peaks are smoothening.

* The optimal bandwidth h can be computed using Silverman's thumb rule . ~~It~~

* It assumes the data is normally distributed.