– So far we have seen parametric densities like Gaussian, GMM, etc,
which makes an assumption about the form.

– non-parametric estimation – estimate $p(x)$ w/o strong assumptions,
using the data. (Note, also has parameters)

## Histogram

- Assume samples $\{x_1, \ldots, x_n\}$
- Consider a region $R$
- Define $P = p(x \in R) = \int_{x \in R} p(x) dx$

  $R$    prob. a point in $R$

- Define $K_R = $ # points inside $R$
- Estimate of $P$:   $\hat{P} = \frac{K_R}{n}$
- Assume $R$ is small, then

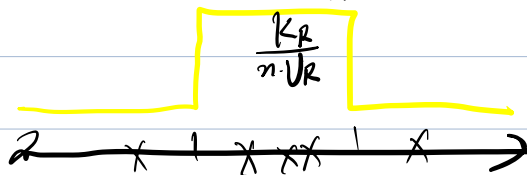$$\hat{P} \approx p(x) \cdot V_R, \quad V_R = \text{volumn of } R,$$
$$x = \text{center of } R.$$

approximate integral over $R$ with rectangle.

$$\underline{\hat{P} = \frac{K_R}{n} \qquad \hat{P} = p(x) \cdot V_R}$$

$$\rightarrow p(x) \cdot V_R = \frac{K_R}{n}$$

$$p(x) = \frac{K_R}{n \cdot V_R}$$

$$\frac{K_R}{n \cdot V_R}$$

This is just a histogram, but we can extend it

Q: How to choose R

✓ 1) keep $V_R$ fixed, & let $K_R$ vary → Parzen windows / Kernal density estimation

2) keep $K_R$ fixed, & let $V_R$ vary → K-NN estimation (k nearest neighbors)

in general 1) is better, why?

Second one requires the same number of points in each region. for a low density region ( region should be extemely large)

dense region ( region will be too small)

Kernal Density Estimation.

- let R be a d-dim hypercube w/ side of h
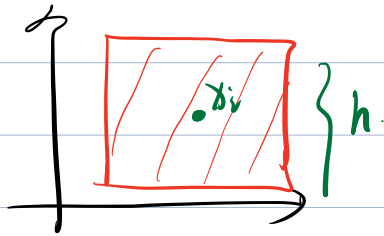
$d=1$ ⟵ h ⟶

$d=2$ ▭ h

$d=3$ ⬡ h

- introduce a window

$$\mathcal{U}(x) = \begin{cases} 1, & |x_i| \leq \frac{1}{2}, \theta i = \{1,...,d\} \\ 0, & \text{otherwise} \end{cases}$$
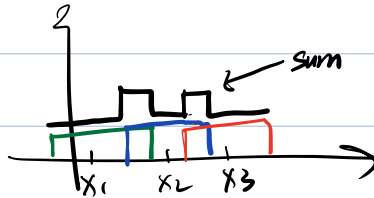
(Parzen window → kernal function)

$$K\left(\frac{x-x_i}{h}\right) = \begin{cases} 1, & \text{if } x \text{ falls inside a cube } w/ \text{ side } h, \\ & \text{centered at } x_i \\ 0, & \text{otherwise} \end{cases}$$
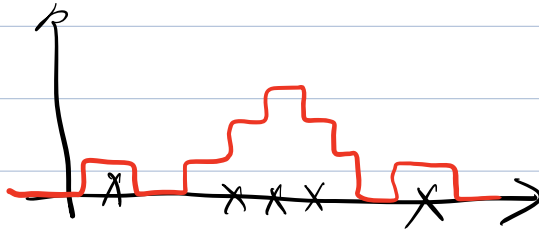
— # of points near $x$: $K = \sum_{i=1}^{n} k\left(\frac{x-x_i}{h}\right)$



sum

Stacking boxes
centered at all $x_i$'s "

$$\hat{p}(x) = \frac{1}{h} \frac{K_R}{V_R}$$

$$\hat{p}(x) = \frac{1}{n \cdot h^d} \sum_{i=1}^{n} k\left(\frac{x-x_i}{h}\right)$$



Other kernel function.

Constraints: $k(x) \geq 0$

$\int k(x)dx = 1$ $\Big\}$ i.e. Valid pdf

Examples:

uniform $k(x) = \begin{cases} 1, & |x_i| \leq \frac{1}{2} \; \forall i \\ 0, & \text{otherwise.} \end{cases}$

unit sphere $k(x) = \begin{cases} \frac{1}{c}, & \|x\|^2 \leq 1 \quad \text{— c is the volume} \\ 0, & 0 \end{cases}$

Gaussian: $k(x) = \frac{1}{\sqrt{2\pi}^d} e^{-\frac{1}{2}\|x\|^2}$   (assume $\sigma^2 = 1$)

$\hat{p}(x) = \frac{1}{n \cdot h^d} \cdot \sum_i k\left(\frac{x - x_i}{h}\right) = \frac{1}{n} \sum_i \mathcal{N}(x | x_i, h^2 I)$
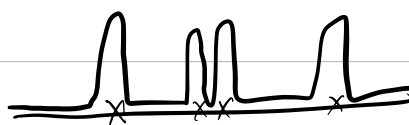
$\frac{1}{n}$ → $\pi_i$

mean, ith comp.

GMM w. $n$ components.

## Bandwidth parameter $h$.

$h$ controls the smoothness of $\hat{p}$
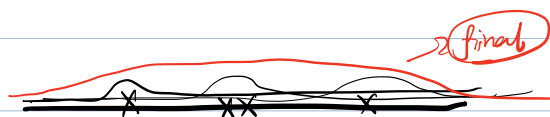
Intuitively

$h$ too small



noisy estimate if
not enough samples

$h$ too large



final

blurry estimate
if too many samples

## Convergence Analysis

Will $\hat{p}(x)$ converge to true $p(x)$?

$\hat{p}(x)$ depends on samples $\{x_i\}_i$, which are r.v.s,

$\Rightarrow$ we can look at bias/variance

$\hat{p}(x)$ converges to $p(x)$ f

① $\lim_{n \to \infty} E[\hat{p}(x)] = p(x)$

② $\lim_{n \to \infty} Var[\hat{p}(x)] = 0$

Define $\tilde{k}(x) = \frac{1}{h^d} k\left(\frac{x}{h}\right) \Rightarrow \hat{p}(x) = \frac{1}{n} \sum_i \tilde{k}(x - x_i)$

Scale the
amplitude

Scale width
of kernal

Mean: $E[\hat{p}(x)] = E_{x_i}\left[\frac{1}{n} \sum_i \tilde{k}(x - x_i)\right]$

⋮             (tutorial

$$= \int p(u) \, \tilde{k}(x-u) \, du$$
$$= p(x) * \tilde{k}(x)$$

<span style="color:red">convolution of the true $p(x)$<br>with the kernel $\tilde{k}(x)$</span>

<span style="color:red">$\Rightarrow$ blurred version of $p(x)$</span>

Only unbiased when

$$\tilde{k}(x) = \delta(x) = \lim_{h \to 0} \tilde{k}(x-x_i) \Rightarrow E[\hat{p}(x)] = p(x)$$

Dirac delta

$$\int f(x) \cdot \delta(x-x_0) \, dx = f(x_0)$$
(formal definition of $\delta$)

Variance : $\text{Var}(\hat{p}(x)) \gtrsim \frac{1}{n h^d} \max_x [k(x)] \cdot E[\hat{p}(x)]$   (tutorial)

For small variance, we need $n$ large / $h$ large.
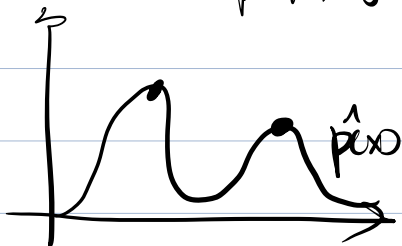
☆  $h$ controls the bias and variance

$\begin{cases} h \to 0 & \text{bias} \to 0 \, , \text{ variance large} \\ h \to \infty & \text{bias} \neq 0 \, , \text{ variance small} \end{cases}$

How to select $h$ ? cross-validation

- select $h$ to maximize $\mathcal{LL}$ of validation set
- select $h$ as function of <u>physical property</u>
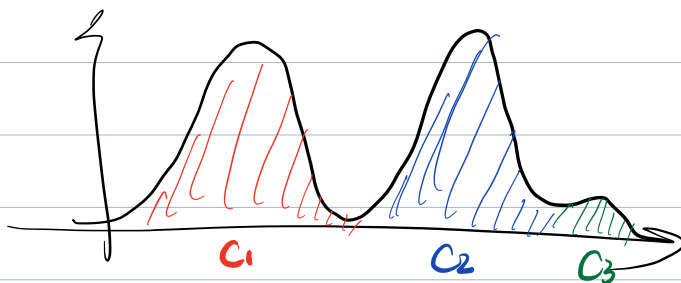    (inherent property of noise)

Mean-Shift algorithm

Find the modes (peaks) of $\hat{p}(x)$


$\hat{p}(x)$

1) Start at a point $\hat{x}$ (eg. one datapoint $x_i$)

2) use gradient ascent to move uphill $[\hat{x} \leftarrow \hat{x} + \lambda \nabla \hat{p}(x)]$

3) eventually $\hat{x}$ will converge to a mode.

• Repeat for many different initial $\hat{x}$s to find the modes

Clustering:

The $x_i$ that converge to the same mode belong to the same cluster.



$C_1$   $C_2$   $C_3$

吸引域的    "basins of attraction"

Consider radially symmetric kernels.

meaning   $\underset{\text{kernel function}}{k(x)} = \underset{\text{Const}}{\alpha} \cdot \underset{\text{kernel profile}}{\overline{k}(\|x\|^2)}$

eg. Gaussian    $k(x) = \frac{1}{(2\pi)^{d/2}} \cdot e^{-\frac{1}{2}\|x\|^2}$

$\overline{k}(r) = e^{-\frac{1}{2}r}$

$\alpha = (2\pi)^{-d/2}$

# KDE (kernel density estimation)

$$\hat{p}(x) = \frac{\alpha}{n \cdot h^d} \sum_i \bar{K}\left(\|\frac{x - x_i}{h}\|^2\right)$$

gradient:  define $\bar{g}(r) = -\bar{K}'(r)$  (Gaussian $\bar{g}(r) = \frac{1}{2} e^{-\frac{1}{2}r}$)

$$\nabla \hat{p}(x) = \left(\frac{\alpha}{n \cdot h^d} \sum_i \bar{K}\left[\left(\frac{x-x_i}{h}\right)^T \left(\frac{x-x_i}{h}\right)\right]\right)'$$

$$= \frac{\alpha}{n \cdot h^{d+2}} \left[\sum_i \bar{g}\left(\|\frac{x-x_i}{h}\|^2\right)(2x - 2x_i)\right]$$

$$= \frac{2\alpha}{n \cdot h^{d+2}} \underbrace{\left(\sum_i \bar{g}\left(\|\frac{x-x_i}{h}\|^2\right)\right)}_{} \cdot \left(\underbrace{\frac{\sum_i x_i \, \bar{g}\left(\|\frac{x-x_i}{h}\|^2\right)}{\sum_i \bar{g}\left(\|\frac{x-x_i}{h}\|^2\right)}}_{} - x\right)$$

consts ↓    kernel profile

large only when $x \to x_i$

$x$ KDE using $\bar{g}(r)$

$$= \hat{g}(x)$$

weighted mean of samples closest to $x$

"mean-shift" vector
diff between weighted mean inside the window and the center of the window.

$$= m(x)$$

Gradient ascent

$$\hat{x}^{(k+1)} = \hat{x}^{(k)} + \lambda \, \nabla \hat{p}(\hat{x}^{(k)})$$

↑ updated    ↑ current    ↑ step size → important for convergence.

Use an adaptive step size

$$\lambda = \frac{1}{\hat{g}(x)} \quad \leftarrow \quad \hat{g}(x) \text{ is small} \Rightarrow \text{large step size}$$

(low-density region)

$$\hat{g}(x) \text{ is large} \Rightarrow \text{small step size}$$

(high-density region)

$$\Rightarrow \quad \hat{x}^{(k+1)} = \hat{x}^{(k)} + \frac{1}{\hat{g}(\hat{x}^{(k)})} \hat{g}(\hat{x}^{(k)}) \cdot m(\hat{x}^{(k)}) \quad \text{(updated gradient descent)}$$

$$= \hat{x}^{(k)} + m(\hat{x}^{(k)})$$

$$= \frac{\sum_i x_i \, \bar{g}\left(\|\frac{\hat{x}^{(k)}-x_i}{h}\|^2\right)}{\sum_i \bar{g}\left(\|\frac{\hat{x}^{(k)}-x_i}{h}\|^2\right)}$$

intuitively



Note: ① proved in paper, guaranteed to converge to

a stationary point of kernal profile $\overline{k}(r)$ is

monotonically decreasing & convex



eg.
Gaussian