

Lecture 5 Non parametric density estimation

So far we have used parametric models - assume a form of the density (Gaussian, Exp, GMM)

Non-parametric estimation - estimate $p(x)$ without assuming a form

Note: non parametric estimation still has parameters

Histogram = samples $\{x_1, \dots, x_N\}$
consider region R

$$\text{Define } \hat{\pi}_R = p(x \in R) = \int_{x \in R} p(x) dx$$

prob: draw sample from R

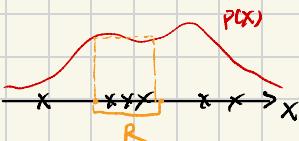
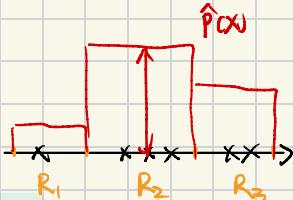
Define: $K_R = \# \text{ samples in } R$

$$MCZ: \hat{\pi}_R = K_R / N$$

Assume R is small ($p(x)$ is flat locally)

$$\hat{\pi}_R = p(x) \cdot V_R \quad \text{volume of } R$$

$$\text{Solve } p(x): p(x) V_R = \frac{K_R}{N} \Rightarrow \hat{p}(x) \approx \frac{K_R}{N \cdot V_R}, x \in R$$



How to Select R ?

- 1) keep V_R fixed, let K_R vary kernel density estimator (KDE)
- 2) keep K_R fixed, let V_R vary k-NN estimator

Kernel Density Estimation

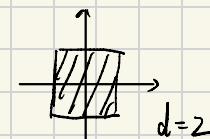
- Hypercube: R d-dim w/ size h



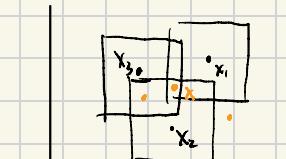
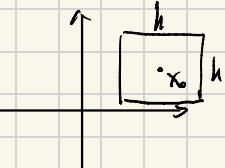
$$V_R = h^d$$

- window function ($h=1$)

$$p(x) = \begin{cases} 1, & |x_i| < \frac{1}{2}, \forall i \in \{1, \dots, d\} \\ 0, & \text{otherwise} \end{cases}$$

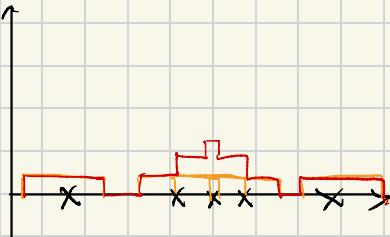


Note: $K\left(\frac{x-x_i}{h}\right) = \begin{cases} 1, & x \text{ is inside hypercube centered at } x_i, \\ & \text{w/ length } h. \\ 0, & \text{otherwise} \end{cases}$



$$K = \sum_{i=1}^N K\left(\frac{x-x_i}{h}\right)$$

$$\text{Hence: } \hat{p}(x) = \frac{1}{N} \frac{K_R}{V_R} = \frac{1}{N h^d} \sum_{i=1}^N K\left(\frac{x-x_i}{h}\right)$$



other kernel func

$$k(x) \geq 0 \\ \int k(x) dx = 1 \quad \} \Rightarrow \text{pdf}$$

Example:

$$\text{unif box } k(x) = \begin{cases} 1 & |x_i| \leq \frac{1}{2}, \forall i \in \{1, \dots, d\} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{unit sphere } k(x) = \begin{cases} \frac{1}{\alpha} & \|x\|^2 \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

α : volume of Hypersphere

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

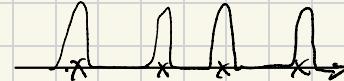
$$\Rightarrow \hat{p}(x) = \frac{1}{N \cdot h^d} \sum K\left(\frac{x - x_i}{h}\right) = \frac{1}{N} \sum N(x | \bar{x}_i, h^2 I)$$

GMM T_{ij} mean, y_j cov Σ_j

bandwidth param

h controls the smoothness of $\hat{p}(x)$

h too small



Noisy estimate if no enough points

h too large



Blurred estimate if the points are too close

Convergence Analysis

- Will $\hat{p}(x)$ converge to the $p(x)$?
 - $\hat{p}(x)$ converges to $p(x)$ if $\lim_{N \rightarrow \infty} \mathbb{E}[\hat{p}(x)] = p(x)$
- 2) $\lim_{N \rightarrow \infty} \text{Var}[\hat{p}(x)] = 0$

Define : $\tilde{k}(x) = \frac{1}{h^d} k\left(\frac{x}{h}\right)$
↑ scale the width
scale amplitude:

$$\Rightarrow \hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \tilde{k}(x - x_i)$$

Mean : $\mathbb{E}[\hat{p}(x)] = \int p(u) \cdot \tilde{k}(x-u) du = p(x) * \tilde{k}(x)$

only unbiased.

$$\tilde{k}(x) = \delta(x) = \lim_{h \rightarrow 0} \tilde{k}(x) \Rightarrow \hat{p}(x) = p(x)$$

Variance

$$\text{var}(\hat{p}(x)) \leq \frac{1}{Nh^d} \left(\max_x k(x) \right) \mathbb{E}[\tilde{k}(x)]$$

Not changeable

For small variance, we need to have h large or N large

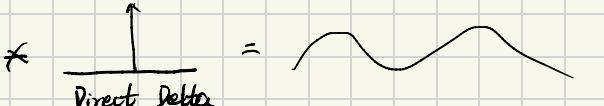
Summary : h controls the trade-off bias vs var.

$$h \rightarrow 0 \Rightarrow \text{bias} = 0, \text{ var} = \infty$$

$$h \rightarrow \infty \Rightarrow \text{bias} = \infty, \text{ var} = 0$$

Convolution between the true $p(x)$ and the kernel

blur the true $p(x)$

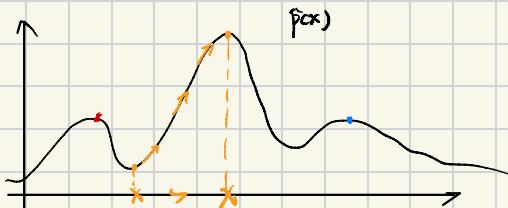


$$\delta(x) = \begin{cases} \infty & , x=0 \\ 0 & , \text{other} \end{cases}$$

$$\int f(x) \delta(x) = f(0)$$

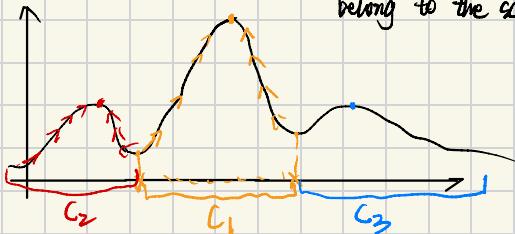
Mean shift algorithm

Find a mode in $\hat{p}(x)$
(peak)



-
- Steps:
- 1) start a random point x_i
 - 2) use grad. ascent to move uphill on $\hat{p}(x)$
 $\hat{x}^{(t+1)} \leftarrow \hat{x}^{(t)} + \lambda \nabla \hat{p}(\hat{x}^{(t)})$
 - 3) repeat, \hat{x} converge to a mode
- repeat over all samples to find all modes
-

clustering: the x_i that share the same mode.
belong to the same cluster



radially symmetric kernel

$$k(x) = c \bar{k}(\|x\|^2)$$

↑
constant ↑
kernel profile distance

Gaussian

$$k(x) = \frac{1}{(2\pi)^d} e^{-\frac{1}{2}\|x\|^2}$$

$$\bar{k}(r) = e^{-\frac{1}{2}r^2}, r = \|x\|^2$$

Density Estimate: $\hat{p}(x) = \frac{c}{N h^d} \sum_{i=1}^N \bar{k}\left(\left\|\frac{x-x_i}{h}\right\|^2\right)$

Gradient: define $\bar{g}(r) = -\bar{k}'(r)$, $\hat{g}(x) = \frac{1}{2} e^{-\frac{1}{2}\|x\|^2}$
(Gaussian)

$$\nabla \hat{p}(x) = \frac{2c}{Nh^d} \cdot \underbrace{\left(\sum_{i=1}^N \bar{g}\left(\left\|\frac{x-x_i}{h}\right\|^2\right) \right)}_{\approx \text{densely estimate using } \bar{g}(r)} \cdot \underbrace{\left(\frac{\sum_i x_i \bar{g}\left(\left\|\frac{x-x_i}{h}\right\|^2\right)}{\sum_i \bar{g}\left(\left\|\frac{x-x_i}{h}\right\|^2\right)} - x \right)}_{\text{point}}$$

nominator: weighted sum
of x_i that are close
to x , weighted by $\bar{g}(\cdot)$

denom: sum of weight

whole: weighted avg
of points close to x

mean shift vector of x , m(x)

gradient ascent

$$\hat{x}^{(k+1)} = \hat{x}^{(k)} + \lambda \cdot \nabla_p (\hat{f}(\hat{x}^{(k)}))$$

$$= \hat{x}^{(k)} + \lambda \cdot \mathbb{L} \cdot \hat{g}(\hat{x}^{(k)}) \cdot m(\hat{x}^{(k)})$$

step size for gradient ascent

adaptive step size

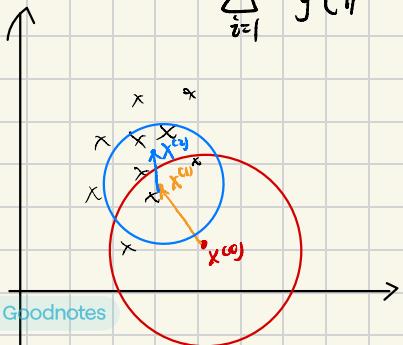
$$\lambda = \frac{1}{\hat{g}(\hat{x}^{(k)}) \cdot \mathbb{L}}$$

*$\hat{g}(x)$ is small \Rightarrow distant
 \Rightarrow large step*
 *$\hat{g}(x)$ is large \Rightarrow close
 \Rightarrow small step*

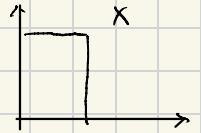
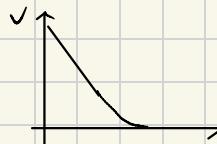
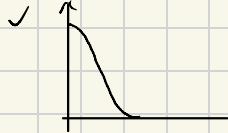
$$\begin{aligned}\hat{x}^{(k+1)} &= \hat{x}^{(k)} + \frac{1}{\hat{g}(\hat{x}^{(k)}) \cdot \mathbb{L}} \cdot \hat{g}(\hat{x}^{(k)}) \cdot m(\hat{x}^{(k)}) \\ &= \hat{x}^{(k)} + m(\hat{x}^{(k)})\end{aligned}$$

$$= \frac{\sum_{i=1}^N x_i \bar{g}(\|\frac{x - x_i}{h}\|^2)}{\sum_{i=1}^N \bar{g}(\|\frac{x - x_i}{h}\|^2)}$$

mean-shift



Note: the algorithms is guaranteed to converge to a stationary point, i.e., find the mode, if the kernel profile is monotonously decrease and convex.



(not convex)