

CSE 547/
STAT 548

4/11/25

Lecture 4

KDE - selecting h
NP clustering
Mean Shift

Project info
hw1 posted

Lecture II – Clustering – Part II: Non-parametric clustering

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

CSE 547/STAT 548
Winter 2022

1 Paradigms for clustering



2 Methods based on non-parametric density estimation



3 Model-based: Dirichlet process mixture models

Kernel density estimation

Input

- data $\mathcal{D} \subseteq \mathbb{R}^d$
- Kernel function $K(z)$
- parameter kernel width h (is a smoothness parameter)

Output

$f(x)$ a probability density over \mathbb{R}^d

• $h = \text{smoothing parameter}$

• $h \uparrow$ Bias↑ Variance↓

fiction call $\rightarrow f(x) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$

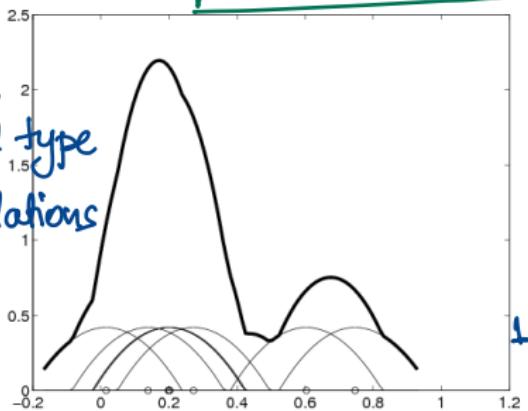
"parameters"

KDE

Train: store $X_{1:n}$
 h , kernel type

"Testing": all calculations

Use



- f is sum of Gaussians centered on each x_i
- f is smoother (less variation) if h larger
- **caveat:** dimension d can't be too large

• $h, K()$ not learned from data

• $\{x_{1:n}\}$ grows with n

h : Model Selection

↳ Choose "grid" of h values

Ex: $\{0.05, 0.25, 0.5, 1, 2, \dots\}$
 logarithmic grid

CV the idea

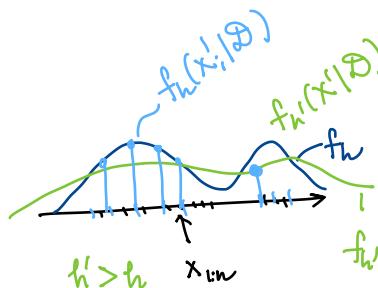
Given $\mathcal{D} = \{x_{1:n}\}$, $K()$
 $H = \{h_1, h_2, \dots, h_K\}$ set of h values ("grid")

1. get sample $\mathcal{D}' = \{x'_{1:n}\} \sim \text{same } P$

2. Compute validation likelihood

for $h \in H$
$$l^V(h) = \ln f_h(\mathcal{D}' | \mathcal{D})$$

3. choose $h^* = \arg \max_{h \in H} l^V(h)$



The kernel function

- Example $K(z) = \frac{1}{(2\pi)^{d/2}} e^{-||z||^2/2}$, $z \in \mathbb{R}^d$ is the Gaussian kernel
- In general
 - $K()$ should represent a density on \mathbb{R}^d , i.e $K(z) \geq 0$ for all z and $\int K(z)dz = 1$
 - $K()$ symmetric around 0, decreasing with $||z||$
- In our case, K must be **differentiable**

I Mean Shift Clustering

(II Level Set Clustering)

I Idea: the mean-shift

$$f(x) = \frac{1}{n h^d} \sum_{i=1}^n e^{-\frac{\|x-x_i\|^2}{2h^2}}$$

$$\nabla f(x) = \frac{1}{n h^d (2\pi)^{d/2}} \sum_{i=1}^n + \frac{2(x-x_i)}{2h^2} \cdot e^{-\frac{\|x-x_i\|^2}{2h^2}} = 0$$

$$\sum_{i=1}^n (x-x_i) K_h(x-x_i) = 0$$

peak

$$\sum_{j=1}^n x_j K_h(x-x_j) = \sum_{i=1}^n x_i K_h(x-x_i)$$

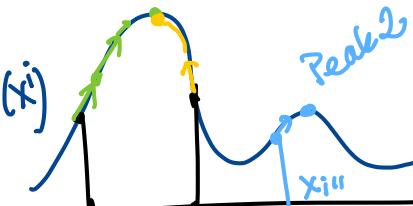
$$x \cdot \left[\sum_{i=1}^n K_h(\) \right] = \sum_{i=1}^n x_i K_h(\)$$

$$x = \sum_{i=1}^n \frac{K_h(x-x_i)}{Z(x)} \cdot x_i = \sum_{i=1}^n w_i x_i \quad \begin{array}{l} \text{weighted sum of } x_{1:n} \\ \text{when } \nabla f(x)=0 \end{array}$$

$w_i(x) > 0$

$$\sum_{i=1}^n w_i = 1$$

Peak(x_i) = Peak 1



h fixed

$$KDE = f(x)$$

Cluster = "data under peak"

$x_i \in \text{Cluster 1}$

$x_{i''} \in \text{Cluster 2}$

Mean-shift

$$m(x) - x = 0 \Leftrightarrow \nabla f(x) = 0$$

Alg.: Iterate $x \leftarrow m(x)$ until convergence

Mean shift algorithms

Idea find points with $\nabla f(x) = 0$

Assume $K(z) = e^{-||z||^2/2}/\sqrt{2\pi}$ Gaussian kernel

$$\nabla f(x) = -\frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)(x-x_i)/h$$

Local max of f is solution of implicit equation

$$x = \underbrace{\frac{\sum_{i=1}^n x_i K\left(\frac{x-x_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)}}_{\text{the mean shift } m(x)}$$

Algorithm Simple Mean Shift

Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, kernel $K(z)$, h

- ① for $i = 1 : n$
 - ① $x \leftarrow x_i$ *initialize at x_i*
 - ② iterate $x \leftarrow m(x)$ until convergence to m_i
- ② group points with same m_i in a cluster *only now known*

Alg properties

1. #peaks K depends on h

2. Runtime

$$\sum_{i=1}^n \frac{k_h(x-x_i)}{z(x)} \cdot x_i$$

$\Theta(n)$

$w_i(x) > 0$

$\sum_{i \in \text{neighbors}(x)}$

\times each step \times all $i = 1:n$
until convergence

\uparrow

$\Theta(n^2)$

1. subsample $n' \ll n$
 \Rightarrow find all peaks $m_{1:k}$
2. assign remaining to nearest m_k

Remarks

- mean shift iteration guaranteed to converge to a max of f
- computationally expensive
- a faster variant...

Algorithm Mean Shift (Comaniciu-Meer)

Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, kernel $K(z)$, h

- ① select q points $\{x_j\}_{j=1:q} = \mathcal{D}_q \subseteq \mathcal{D}$
that cover the data well
- ② for $j \in \mathcal{D}_q$
 - ① $x \leftarrow x_j$
 - ② iterate $x \leftarrow m(x)$ until convergence to m_j
- ③ group points in \mathcal{D}_q with same m_j in a cluster
- ④ assign points in $\mathcal{D} \setminus \mathcal{D}_q$ to the clusters by the **nearest-neighbor** method

$$k(i) = k(\operatorname{argmin}_{j \in \mathcal{D}_q} \|x_i - x_j\|)$$