CSE 547 STAT 548

4/14/25



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

Marina Meila (UW)



DBSCAN

 < □ > < ⊡ > < ⊡ > < ⊡ > < ⊡ > < ⊡ > < ⊡ > < ⊡ > < ⊡ > < ⊡ > < ○ < ○</td>

 CSE 547/STAT 548 Winter 2022
 2/

Marina Meila (UW)

What is clustering? Problem and Notation

- Informal definition Clustering = Finding groups in data
- Notation $\mathcal{D} = \{x_1, x_2, \dots x_n\}$ a data set
 - *n* = number of **data points**
 - K = number of clusters ($K \ll n$)
 - $\Delta = \{C_1, C_2, \dots, C_K\} \text{ a partition of } \mathcal{D} \text{ into disjoint subsets}$
 - k(i) = the label of point *i*
 - $\mathcal{L}(\Delta) = \text{cost (loss) of } \Delta \text{ (to be minimized)}$
- Second informal definition Clustering = given *n* data points, separate them into *K* clusters
- Hard vs. soft clusterings
 - \bullet Hard clustering $\Delta:$ an item belongs to only 1 cluster
 - Soft clustering $\gamma = \{\gamma_{ki}\}_{k=1:K}^{i=1:n}$

 γ_{ki} = the degree of membership of point *i* to cluster *k*

$$\sum_k \gamma_{ki} = 1 \quad \text{for all } i$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

CSE 547/STAT 548 Winter 2022

3/

(usually associated with a probabilistic model)

Methods based on non-parametric density estimation

Idea The clusters are the isolated peaks in the (empirical) data density

イロト 不得 トイヨト イヨト

CSE 547/STAT 548 Winter 2022

э

6/

- group points by the peak they are under
- some outliers possible
- *K* = 1 possible(no clusters)
- shape and number of clusters K determined by algorithm
- structural parameters
 - smoothness of the density estimate
 - what is a peak

Algorithms

- peak finding algorithms Mean-shift algorithms
- level sets based algorithms
 - Nugent-Stuetzle, Support Vector clustering
- Information Bottleneck ?

Kernel density estimation

Input

data D ⊆ ℝ^d
 Kernel function K(z)

- parameter kernel width h (is a smoothness parameter)
- **Putput** f(x) a **probability density** over \mathbb{R}^d

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



< □ > < □ > < □ > < □ > < □ >

CSE 547/STAT 548 Winter 2022

7/

- 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

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 $K(z) \ge 0$ for all z and $\int K(z)dz = 1$
 - K() symmetric around 0, decreasing with ||z||
- In our case, K must be differentiable

Mean shift algorithms

Idea find points with $\nabla f(x) = 0$ Assume $K(z) = e^{-||z||^2/2}/\sqrt{2\pi}$ Gaussian kernel

$$abla f(x) = -rac{1}{nh^d}\sum_{i=1}^n \mathcal{K}(rac{x-x_i}{h})(\mathbf{x}-x_i)/h$$

Local max of f is solution of implicit equation

$$x = \underbrace{\frac{\sum_{i=1}^{n} x_i K(\frac{x-x_i}{h})}{\sum_{i=1}^{n} K(\frac{x-x_i}{h})}}$$

the mean $shift_{m(x)}$

◆□ ▶ ◆□ ▶ ◆ 三 ▶ ◆ 三 ▶ ● ○ ○ ○ ○ ○ CSE 547/STAT 548 Winter 2022

9/

Algorithm Simple Mean Shift Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, kernel K(z), h • for $i = 1 : \tilde{n}$ $() x \leftarrow x_i$ 2 iterate $x \leftarrow m(x)$ until convergence to m_i **2** group points with same m_i in a cluster

. MeanSh $(x_i) \rightarrow \mu_e \Leftrightarrow X_i \in C_e$





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

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

◆□▶ ◆□▶ ◆三▶ ◆三▶ ・三 ・ のへで CSE 547/STAT 548 Winter 2022

[Gaussian blurring mean shift]

Idea

• like Simple Mean Shift but points are shifted to new locations

イロン イ団 とく ヨン イヨン

CSE 547/STAT 548 Winter 2022

э

11

- the density estimate f changes
- becomes concentrated around peaks very fast

Algorithm Gaussian Blurrring Mean Shift (GBMS)

- Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, Gaussian kernel K(z), h
 - Iterate until STOP
 - for i = 1 : n compute $m(x_i)$
 - $o for i = 1 : n, x_i \leftarrow m(x_i)$

Remarks

- all x_i converge to a single point
 - \Rightarrow need to stop before convergence

Empirical stopping criterion ?

- define $e_i^t = ||x_i^t x_i^{t-1}||$ the change in x_i at t
- define $H(e^t)$ the entropy of the histogram of $\{e_i^t\}$
- STOP when $\sum_{i=1}^{n} e_i^t / n < tol OR |H(e^t) H(e^{t-1})| < tol'$

Convergence rate If true f Gaussian, convergence is cubic

$$||x_i^t - x^*|| \le C ||x_i^{t-1} - x^*||^3$$

very fast!!







The Nugent-Stuetzle algorithm

Algorithm Nugent-Stuetzle

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

- Compute KDE f(x) for chosen h
- (a) for levels $0 < l_1 < l_2 < \ldots < l_r < \ldots < l_R \ge \sup_x f(x)$
 - find level set $L_r = \{x \mid f(x) \ge l_r\}$ of f

Q if L_r disconnected then each connected component is a cluster $\rightarrow (C_{r,1}, C_{r,2}, \dots, C_{r,\kappa_r})$

Solutput clusters $\{(C_{r,1}, C_{r,2}, \dots, C_{r,K_r})\}_{r=1:R}$

Remarks

- every cluster $C_{r,k} \subseteq$ some cluster $C_{r-1,k'}$
- therefore output is hierarchical clustering
- some levels can be pruned (if no change, i.e. $K_r = K_{r-1}$)

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

э

- algorithm can be made recursive, i.e. efficient
- finding level sets of f tractable only for d = 1, 2
- for larger d, $L_r = \{x_i \in \mathcal{D} \mid f(x_i) \geq I_r\}$
- to find connected components
 - for $i \neq j \in L_r$ if $f(tx_i + (1 - t)x_j) \ge l_r$ for $t \in [0, 1]$ then k(i) = k(j)
- confidence intervals possile by resampling

Cluster tree with 13 leaves (8 clusters, 5 artifacts)



メロト メタト メヨト メヨト

CSE 547/STAT 548 Winter 2022

æ

15

(from ?)

Chaudhuri-Dasgupta Algorithm



The K-nn density estimator

The K-nn density estimator

- Let $B_r(x)$ be the (closed) ball of radius r centered at x
- If $|B_r(x^i) \cap D| = k$ then $\hat{p}(x^i) = \frac{1}{r^n \omega_n} \frac{k}{n}$ is an estimate of the density at x^i
 - $\omega_n = \pi^{n/2} / \Gamma(n/2 + 1)$ is the volume of the unit ball in \mathbb{R}^n
 - intuitively, the ball of radius r contains k/n probability mass
 - Note that the density \hat{p} is not required to integrate to 1

p(x) = K-nn density est. le=parameter smoothing)

Simple -> good for large n DBScan



- Algorithm idea
- Construct directed graph \mathcal{G} with edges (i, j) where $x^i \in Q, j \in B_r(x^i)$
- The graph edges between core points are undirected/symmetric, the other are from core to border
- Clusters are determined by the connected components of the graph restricted to Q.
- The border points are assigned to a cluster containing x^j so that $x^i \in B_r(x^j), x^j \in Q$ Note a. that this assignment is not unique!
 - Heuristic algorithm estimates r, m

Marina Meila (UW)

[Chaudhuri-Dasgupta Algorithm]

Consistency Theorem For any ϵ (separation parameter) and δ (confidence), $\alpha \in [\sqrt{2}, 2]$ (graph density), if $k = C \log^2(1/\delta) \frac{d\log n}{\epsilon^2}$

for any two clusters C, C' in cluster tree, there exists a level r so that $C \cap D, C' \cap D$ are clusters at level r

• r depends on $\lambda =$ "bridge" between C, C' (and $\sigma > 0$ "tube" width)

$$r^d \omega_d \lambda = \frac{k}{n} + \dots$$
 confidence term

• it follows that the needed sample size n at level λ

$$n \,=\, \mathcal{O}\left(rac{d}{\lambda\epsilon^2(\sigma/2)^d\omega_d}\lograc{d}{\lambda\epsilon^2(\sigma/2)^d\omega_d}
ight)$$

イロト 不得 トイヨト イヨト 三日

CSE 547/STAT 548 Winter 2022

- this sample complexity n is almost tight
- for $\alpha < \sqrt{2}$ sample complexity is exponential in d
- New results [Kent, B. P., Rinaldo, A. and Verstynen, T. 2013]
- Remark: algorithm(s) can be applied in any metric space