

BIOSTAT527

4/12/23

# Lecture 6

Non-parametric  
clustering

# Lecture IV – Non-parametric clustering

Marina Meilă  
[mmp@stat.washington.edu](mailto:mmp@stat.washington.edu)

Department of Statistics  
University of Washington

soft  
hard

vectors  
graphs

param (k-means, EM)

NP — density estimation

① Paradigms for clustering



② Methods based on non-parametric density estimation



③ Model-based: [ Dirichlet process mixture models ]

Reading AoNPS Ch.: –, HTF Ch.: 14.3 Murphy Ch.: 11.[1], 11.2.1-3, 11.3, Ch 25

# Clustering Paradigms

Depend on type of data, type of clustering, type of cost (probabilistic or not), and constraints (about  $K$ , shape of clusters)

- Data = vectors  $\{x_i\}$  in  $\mathbb{R}^d$

Parametric ( $K$ known)	Cost based [hard]
	Model based [soft]

Non-parametric ( $K$ determined by algorithm)	Dirichlet process mixtures [soft]
	Information bottleneck [soft]
	Modes of distribution [hard]
	Gaussian blurring mean shift? [hard]
	Level sets of distribution [hard]

- Data = similarities between pairs of points  $[S_{ij}]_{i,j=1:n}$ ,  $S_{ij} = S_{ji} \geq 0$

## Similarity based clustering

Graph partitioning	spectral clustering [hard, $K$ fixed, cost based]
	typical cuts [hard non-parametric, cost based]
Affinity propagation	[hard/soft non-parametric]

# Methods based on non-parametric density estimation

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

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

$K \text{ SDE}$

## Mean shift algorithms

$$f(x) = \text{KDE}_h$$

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

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

$h$

$b(z)$

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

$$k'(z) = -z k(z) \quad z \in \mathbb{R}^d$$

Local max of  $f$  is solution of implicit equation

Mean Shift

$m(x) - x$

$$x = \frac{\sum_{i=1}^n x_i K\left(\frac{x-x_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)} = \sum_i w_i x^i = m(x)$$

the mean shift  $m(x)$

weighted mean around  $x$  = a peak  $\mu_1$

Algorithm Simple Mean Shift

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

① for  $i = 1 : n$

    ①  $x \leftarrow x_i$

    ② iterate  $x \leftarrow m(x)$  until convergence to  $m_i$

② group points with same  $m_i$  in a cluster



$$\text{Ex: } m(x) - x \propto \nabla f(x)$$

## Computation

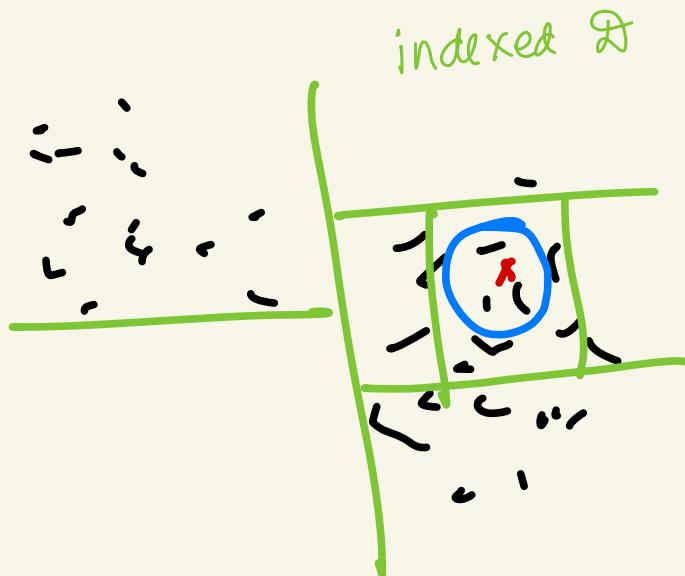
$m(x)$   $n \cdot k()$  evaluations

$| \{ \text{neighbors of } x \} | \cdot k()$  evaluations

if method to find neighbors used

$\times$  iterations to local max  $\mu_k$

$x \cdot n \cdot \mu_k$  for every  $x^i \in D$



$$\nabla f(x) = 0$$

$$x \cdot \sum k\left(\frac{x-x^i}{h}\right) = \sum k\left(\frac{x-x^i}{h}\right) \cdot x^i$$

KDE  
 $\approx f(x)$

weighted avg of  $x^i$  around  $x$

$$x = \frac{\sum_i k(\cdot) x^i}{\sum_i k(\cdot)}$$

$w_i \equiv \text{Nadaraya-Watson}$

## 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  $K(z)$ ,  $h$ 

heuristic acceleration

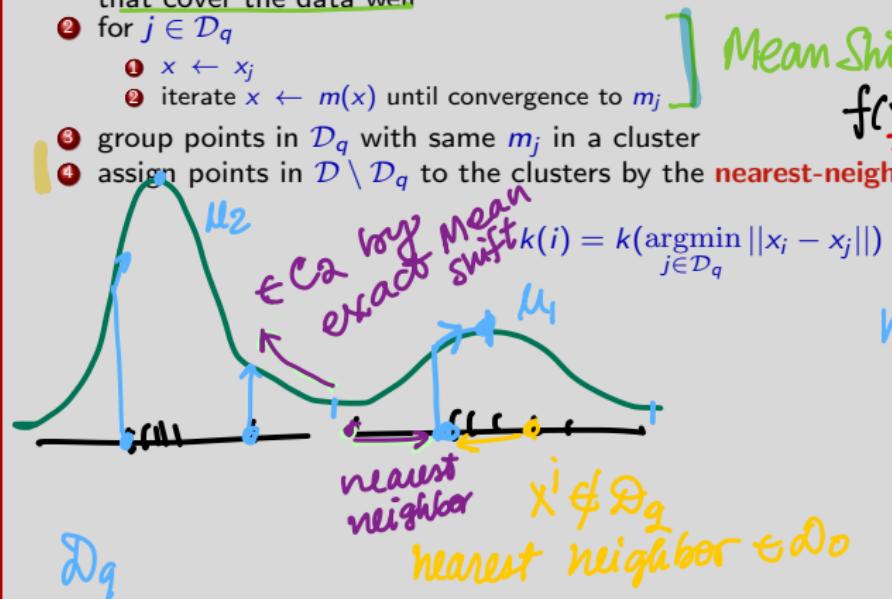
① select  $q$  points  $\{x_j\}_{j=1:q} = \mathcal{D}_q \subseteq \mathcal{D}$  ← landmarks  
that cover the data well

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

Mean Shift

 $f(x)$  from all data

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

$$n_g = |\mathcal{D}_g|$$

$n_g \times n \times$  iteration  
↑ KBS

$$(n - n_g) \times n_g$$

## Farthest points heuristic (Anchor algorithm)

Given  $\mathcal{D} = \{x^1, \dots, x^n\}$

Want  $\mathcal{D}' \subset \mathcal{D}$  "that covers  $\mathcal{D}$  well"

$$n' = |\mathcal{D}'| \quad \mu^{1:n'}$$

$\mu_1 \sim \text{uniform}(\mathcal{D})$

for  $j = 1, 2, \dots, n'$

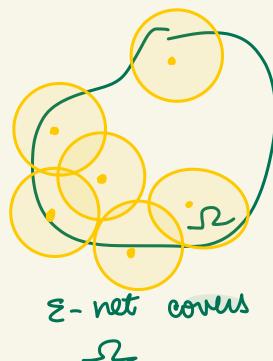
$$\mu_j = \underset{x^i \notin \mathcal{D}'}{\operatorname{argmax}} \left( \min_{i'=1:j-1} \|x^i - \mu_{i'}\| \right)$$

furthest closest neighbor in  $\mathcal{D}'$

### Problems

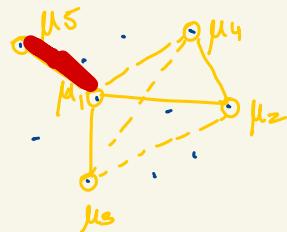
- selects outliers !!

-  $n' = ?$



$\Rightarrow \forall x \in \mathcal{D}$  has

$$\|x - \mu_i\| < \varepsilon \quad \text{for some } \mu_i$$



## [Supplement: Gaussian blurring mean shift]

### Idea

- like **Simple Mean Shift** but points are shifted to new locations
- the density estimate  $f$  changes
- becomes concentrated around peaks very fast

### Algorithm Gaussian Blurring Mean Shift (GBMS)

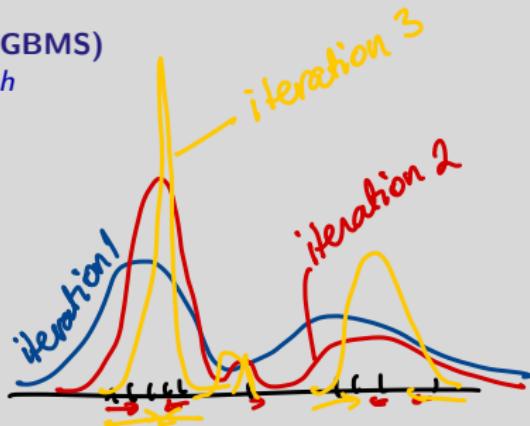
**Input** Data  $\mathcal{D} = \{x_i\}_{i=1:n}$ , Gaussian kernel  $K(z)$ ,  $h$

① Iterate until **STOP**

- ① for  $i = 1 : n$  compute  $m(x_i)$
- ② 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



- very fast
- heuristic
- clustering alg

## 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  $\sum_{i=1}^n e_i^t / n < \text{tol}$  OR  $|H(e^t) - H(e^{t-1})| < \text{tol}$ ,

**Convergence rate** If true  $f$  Gaussian, convergence is **cubic**

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

very fast!!

# The Nugent-Stuetzle algorithm

QW!

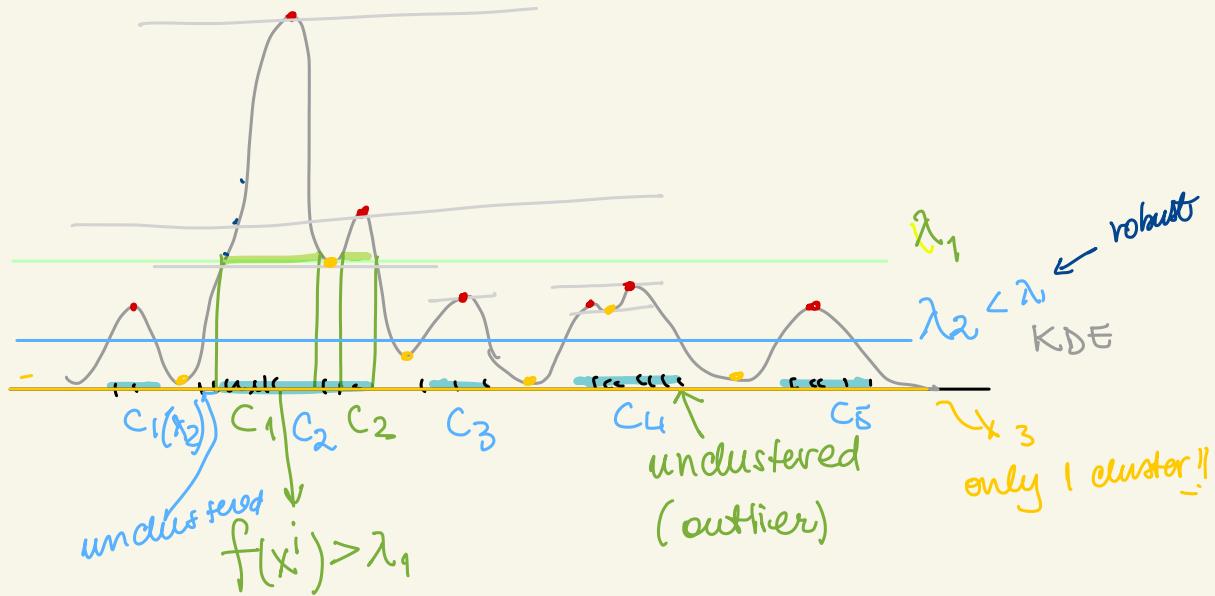
## Algorithm Nugent-Stuetzle

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

- ① Compute KDE  $f(x)$  for chosen  $h$
- ② for levels  $0 < l_1 < l_2 < \dots < l_r < \dots < l_R \geq \sup_x f(x)$ 
  - ① find level set  $L_r = \{x \mid f(x) \geq l_r\}$  of  $f$
  - ② if  $L_r$  disconnected then each connected component is a cluster  $\rightarrow (C_{r,1}, C_{r,2}, \dots, C_{r,K_r})$

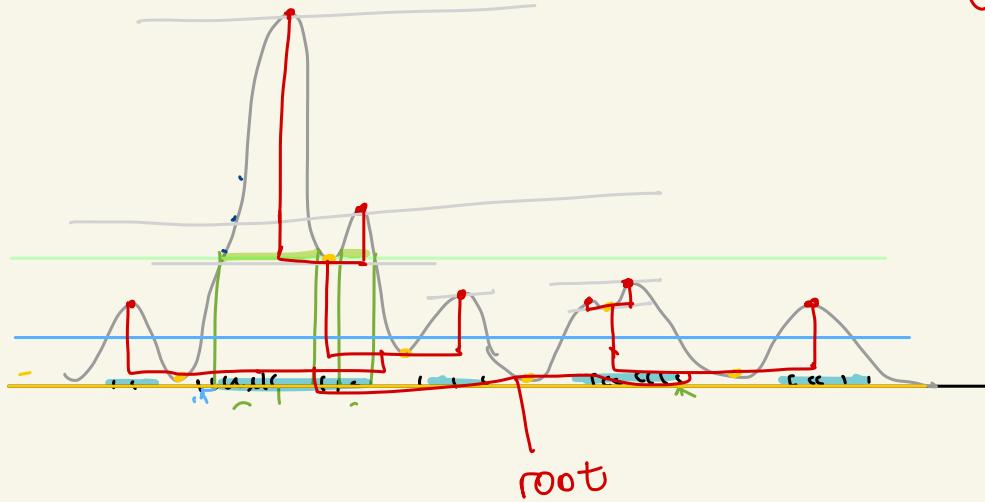
**Output** clusters  $\{(C_{r,1}, C_{r,2}, \dots, C_{r,K_r})\}_{r=1:R}$

1. Compute  $f(x^i)$   $i=1:n$   $n^2$   
 2. ...



- How to choose  $\lambda$  ? - not too many undesired  
- robust  $\lambda \pm \Delta\lambda = \text{same clusters}$

Cluster tree



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