### Lecture II - Clustering - Part I: Parametric clustering

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

CSE 547/STAT 548 Winter 2022

Marina Meila (UW)



#### Paradigms for clustering

Parametric clustering algorithms (K given)

- Cost based / hard clustering
- Model based / soft clustering
- Outliers

Reading MMDS Ch.: 7.3 K-means HTF Ch.:14.3, Murphy Ch.: 11.[1], 11.2.1-3, 11.3, Ch 25

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

3

2/

### 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 ( $K \ll n$ )
  - $\Delta = \{C_1, C_2, \dots, C_K\}$  a partition of  $\mathcal{D}$  into disjoint subsets
  - k(i) = the label of point *i*
  - $L(\Delta) = \cos(\cos) \operatorname{of} \Delta$  (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  $\gamma = \{\gamma_{ki}\}_{k=1:K}^{i=1:n}$

 $\gamma_{ki}$  = the degree of membership of point *i* to cluster *k* 

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

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

3

3/

(usually associated with a probabilistic model)



### 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     | Cost based [hard]           |
|   | (K known)      | Model based [soft           |

| Non-parametric        | Dirichlet process mixtures [soft]                           |
|-----------------------|---|
| ( <i>K</i> determined | Information bottleneck [soft]                               |
| by algorithm)         | Modes of distribution [hard]                                |
|                       | Gaussian blurring mean shiftCarreira-Perpinan (2007) [hard] |

Data = similarities between pairs of points [S<sub>ij</sub>]<sub>i,j=1:n</sub>, S<sub>ij</sub> = S<sub>ji</sub> ≥ 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]                      |

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

## Classification vs Clustering

|                  | Classification       | Clustering                           |
|------------------|----------------------|--------------------------------------|
| Cost (or Loss) L | Expectd error        | many! (probabilistic or not)         |
|                  | Supervised           | Unsupervised                         |
| Generalization   | Performance on new   | Performance on current               |
|                  | data is what matters | data is what matters                 |
| K                | Known                | Unknown                              |
| "Goal"           | Prediction           | Exploration Lots of data to explore! |
| Stage            | Mature               | Still young                          |
| of field         |                      |                                      |

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

э

7/

## Parametric clustering algorithms

- Cost based
  - Single linkage (min spanning tree)
  - Min diameter
    - Fastest first traversal (HS initialization)
  - K-medians
  - K-means
- Model based (cost is derived from likelihood)
  - EM algorithm
  - "Computer science" /" Probably correct" algorithms

## Minimum diameter clustering

• Cost 
$$L(\Delta) = \max_{k} \max_{\substack{i,j \in C_k}} ||x_i - x_j||$$

diameter

- · Mimimize the diameter of the clusters
- Optimizing this cost is NP-hard
- Algorithms

• Fastest First Traversal Hochbaum and Shmoys (1985) – a factor 2 approximation for the min cost For every  $\mathcal{D}$ , FFT produces a  $\Delta$  so that

$$L^{opt} \leq L(\Delta) \leq 2L^{opt}$$

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

э

8/

rediscovered many times

#### **Algorithm Fastest First Traversal**

Input Data  $\mathcal{D} = \{x_i\}_{i=1:n}$ , number clusters K defines centers  $\mu_{1:K} \in \mathcal{D}$ 

(many other clustering algorithms use centers)

- **()** pick  $\mu_1$  at random from  $\mathcal{D}$
- **(a)** for k = 2 : K

 $\mu_k \leftarrow \operatorname{argmax}_{\mathcal{D}} \operatorname{distance}(x_i, \{\mu_{1:k-1}\})$ 

• for i = 1 : n (assign points to centers)

k(i) = k if  $\mu_k$  is the nearest center to  $x_i$ 

## Model based clustering: Mixture models



#### • The mixture density

$$f(x) = \sum_{k=1}^{K} \pi_k f_k(x)$$

•  $f_k(x)$  = the **components** of the mixture

- · each is a density
- f called mixture of Gaussians if  $f_k = Normal_{\mu_k, \Sigma_k}$

< □ > < □ > < □ > < □ > < □ >

CSE 547/STAT 548 Winter 2022

•  $\pi_k$  = the mixing proportions,

$$\sum_k = 1^K \pi_k = 1, \ \pi_k \ge 0.$$

• model parameters  $\theta = (\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K})$ 

## Model based clustering: Mixture models



#### • The mixture density

$$f(x) = \sum_{k=1}^{K} \pi_k f_k(x)$$

•  $f_k(x)$  = the **components** of the mixture

- · each is a density
- f called mixture of Gaussians if  $f_k = Normal_{\mu_k, \Sigma_k}$
- $\pi_k$  = the mixing proportions,  $\sum_k = 1^K \pi_k = 1, \ \pi_k \ge 0.$
- model parameters  $\theta = (\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K})$
- The degree of membership of point *i* to cluster *k*

$$\gamma_{ki} \stackrel{\text{def}}{=} P[x_i \in C_k] = \frac{\pi_k f_k(x)}{f(x)} \text{ for } i = 1: n, k = 1: K$$
(1)

• depends on  $x_i$  and on the model parameters

イロト イボト イヨト イヨト

## The Maximum Likelihood Principle

- Given data  $\mathcal{D} = \{x_{1:n}\}$  sampled i.i.d from some unknown  $P^*$
- Model  $P_{\theta}(x)$  depends on parameter  $\theta$
- Problem: How to estimate  $\theta$ ?
- Principle: Maximum Likelihood

$$\mathsf{Likelihood}(\theta|\mathcal{D}) = P_{\theta}(\mathcal{D}) = \prod_{i=1}^{n} P_{\theta}(x_i)$$

• Often convenient to use log-likelihood  $I(\theta)$ 

$$I(\theta) = \sum_{i=1}^{n} \ln P_{\theta}(x_i)$$

• Reason: many  $P_{\theta}$  are expressed with exponential functions (e.g the Normal distribution)

## Criterion for clustering: Max likelihood

- denote  $\theta = (\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K})$  (the parameters of the mixture model)
- Define likelihood  $P[\mathcal{D}|\theta] = \prod_{i=1}^{n} f(x_i)$
- Typically, we use the log likelihood

$$I(\theta) = \ln \prod_{i=1}^{n} f(x_i) = \sum_{i=1}^{n} \ln \sum_{k} \pi_k f_k(x_i)$$
(2)

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

3

- denote  $\theta^{ML} = \operatorname{argmax}_{\theta} I(\theta)$
- $\theta^{ML}$  determines a soft clustering  $\gamma$  by (1)
- a soft clustering  $\gamma$  determines a  $\theta$  (see later)
- Therefore we can write

$$L(\gamma) = -I(\theta(\gamma))$$

## Algorithms for model-based clustering

Maximize the (log-)likelihood w.r.t  $\theta$ 

- directly (e.g by gradient ascent in  $\theta$ )
- by the EM algorithm (very popular!)
- indirectly, w.h.p. by "computer science" algorithms

w.h.p = with high probability (over data sets)

## The Expectation-Maximization (EM) Algorithm

#### Algorithm Expectation-Maximization (EM)

**Input** Data  $\mathcal{D} = \{x_i\}_{i=1:n}$ , number clusters *K* tialize parameters  $\pi_{1:K} \in \mathbb{R}, \ \mu_{1:K} \in \mathbb{R}^d, \ \Sigma_{1:K} \in \mathbb{R}^{d \times d}$  at random<sup>1</sup> terate until convergence

**E step** (Optimize clustering) for i = 1 : n, k = 1 : K

$$\gamma_{ki} = \frac{\pi_k f_k(x)}{f(x)}$$

**M** step (Optimize parameters) set  $\Gamma_k = \sum_{i=1}^n \gamma_{ki}$ , k = 1 : K (number of points in cluster k)

$$\pi_{k} = \frac{\Gamma_{k}}{n}, \quad k = 1 : K$$

$$\mu_{k} = \sum_{i=1}^{n} \frac{\gamma_{ki}}{\Gamma_{k}} x_{i}$$

$$\Sigma_{k} = \frac{\sum_{i=1}^{n} \gamma_{ki} (x_{i} - \mu_{k}) (x_{i} - \mu_{k})^{T}}{\Gamma_{k}}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで CSE 547/STAT 548 Winter 2022

14

•  $\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K}$  are the maximizers of  $l_c(\theta)$  in (6) •  $\sum_{k} \Gamma_{k} = n$ 

 $<sup>{}^{1}\</sup>Sigma_{k}$  need to be symmetric, positive definite matrices

### The EM Algorithm – Motivation

• Define the indicator variables

$$z_{ik} = \begin{cases} 1 & \text{if } i \in C_k \\ 0 & \text{if } i \notin C_k \end{cases}$$
(3)

denote  $\bar{z} = \{z_{ki}\}_{k=1:K}^{i=1:n}$ • Define the complete log-likelihood

$$l_{c}(\theta, \bar{z}) = \sum_{i=1}^{n} \sum_{k=1}^{K} z_{ki} \ln \pi_{k} f_{k}(x_{i})$$
(4)

• 
$$E[z_{ki}] = \gamma_{ki}$$
  
• Then

$$E[l_{c}(\theta, \bar{z})] = \sum_{i=1}^{n} \sum_{k=1}^{K} E[z_{ki}] [\ln \pi_{k} + \ln f_{k}(x_{i})]$$
(5)

$$= \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_{ki} \ln \pi_{k} + \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_{ki} \ln f_{k}(x_{i})]$$
(6)

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

E 99€

- If  $\theta$  known,  $\gamma_{ki}$  can be obtained by (1) (Expectation)
- If  $\gamma_{ki}$  known,  $\pi_k, \mu_k, \Sigma_k$  can be obtained by separately maximizing the terms of  $E[I_c]$  (Maximization)

イロト 不得 トイヨト イヨト

CSE 547/STAT 548 Winter 2022

э

17

## Brief analysis of EM

$$Q(\theta,\gamma) = \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_{ki} \ln \underbrace{\pi_k f_k(x_i)}_{\theta}$$

- each step of EM increases  $Q(\theta, \gamma)$
- Q converges to a local maximum
- $\bullet$  at every local maxi of  $\textit{Q},\,\theta\,\leftrightarrow\,\gamma$  are fixed point
- $Q(\theta^*,\gamma^*)$  local max for  $Q \Rightarrow I(\theta^*)$  local max for  $I(\theta)$
- under certain regularity conditions  $\theta \longrightarrow \theta^{ML}$  McLachlan and Krishnan (1997)
- the E and M steps can be seen as projections Neal and Hinton (1998)
- Exact maximization in M step is not essential. Sufficient to increase Q. This is called Generalized EM

イロト 不得 トイヨト イヨト

CSE 547/STAT 548 Winter 2022

-

18

### Probablistic alternate projection view of EMNeal and Hinton (1998)

- let  $z_i$  = which gaussian generated *i*? (random variable),  $X = (x_{1:n})$ ,  $Z = (z_{1:n})$
- Redefine Q

$$Q(\tilde{P}, \theta) = L(\theta) - KL(\tilde{P}||P(Z|X, \theta))$$

where  $P(X, Z|\theta) = \prod_i \prod_k P[z_i = k]P[x_i|\theta_k]$  $\tilde{P}(Z)$  is any distribution over Z,  $KL(P(w)||Q(w)) = \sum_w P(w) \ln \frac{P(w)}{Q(w)}$  the Kullbach-Leibler divergence

Then,

- E step  $\max_{\tilde{P}} Q \Leftrightarrow KL(\tilde{P}||P(Z|X,\theta))$
- M step  $\max_{\theta} Q \Leftrightarrow KL(P(X|Z, \theta^{old})||P(X|\theta))$
- Interpretation: KL is "distance", "shortest distance" = projection

### The M step in special cases

• Note that the expressions for  $\mu_k, \Sigma_k$  = expressions for  $\mu, \Sigma$  in the normal distribution, with data points  $x_i$  weighted by  $\frac{\gamma_{ki}}{\Gamma_k}$ 

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

19

|   | M step   |
|---|--|
| general case  | $\Sigma_k = \sum_{i=1}^n \frac{\gamma_{ki}}{\Gamma_k} (x_i - \mu_k) (x_i - \mu_k)^T$   |
| $\Sigma_k = \Sigma$<br>"same shape & size" clusters | $\boldsymbol{\Sigma} \leftarrow \frac{\sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_{ki} (x_{i} - \mu_{k}) (x_{i} - \mu_{k})^{T}}{n}$ |
| $\Sigma_k = \sigma_k^2 I_d$                         | $\sigma_k^2 \leftarrow \frac{\sum_{i=1}^n \gamma_{ki}   \mathbf{x}_i - \mu_k  ^2}{d\Gamma_k}$                                |
| "round" clusters                                    |  |
| $\sum_{k} = \sigma^2 I_d$                           | $\sigma^2 \leftarrow \frac{\sum_{i=1}^n \sum_{k=1}^K \gamma_{ki}   x_i - \mu_k  ^2}{nd}$                                     |
| round, same size clusters                           |  |

Exercise Prove the formulas above

• Note also that K-means is EM with  $\Sigma_k = \sigma^2 I_d, \ \sigma^2 \to 0$  Exercise Prove it



More special cases Banfield and Raftery (1993) introduce the following description for a covariance matrice in terms of volume, shape, alignment with axes (=determinant, trace, e-vectors). The letters below mean: l=unitary (shape, axes), E=equal (for all k), V=unequal

- EII: equal volume, round shape (spherical covariance)
- VII: varying volume, round shape (spherical covariance)
- EEI: equal volume, equal shape, axis parallel orientation (diagonal covariance)
- VEI: varying volume, equal shape, axis parallel orientation (diagonal covariance)
- EVI: equal volume, varying shape, axis parallel orientation (diagonal covariance)
- VVI: varying volume, varying shape, equal orientation (diagonal covariance)
- EEE: equal volume, equal shape, equal orientation (ellipsoidal covariance)
- EEV: equal volume, equal shape, varying orientation (ellipsoidal covariance)
- VEV: varying volume, equal shape, varying orientation (ellipsoidal covariance)
- VVV: varying volume, varying shape, varying orientation (ellipsoidal covariance)

(from Nugent and Meila (2010))

Marina Meila (UW)

20

イロン イ団 とく ヨン イヨン

#### EM versus K-means

- Alternates between cluster assignments and parameter estimation
- Cluster assignments  $\gamma_{ki}$  are probabilistic
- Cluster parametrization more flexible



• Converges to local optimum of log-likelihood Initialization recommended by K-logK method

• Modern algorithms with guarantees (for e.g. mixtures of Gaussians)

- Random projections
- Projection on principal subspace Vempala and Wang (2004)
- Two step EM (=K-logK initialization + one more EM iteration)

< □ > < 同 > < 回 > < 回 >

CSE 547/STAT 548 Winter 2022

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

# "Computer science" algorithms for mixture models

- Assume clusters well-separated
  - e.g  $||\mu_k \mu_l|| \geq C \max(\sigma_k, \sigma_l)$
  - with  $\sigma_k^2 = \max \operatorname{eigenvalue}(\Sigma_k)$
- true distribution is mixture
  - of Gaussians
  - of log-concave  $f_k$ 's (i.e. In  $f_k$  is concave function)
- then, w.h.p. (n, K, d, C)
  - we can label all data points correctly
  - $\bullet \ \Rightarrow \ {\rm we \ can \ find \ good \ estimate \ for \ } \theta$

Even with (S) this is not an easy task in high dimensions

Because  $f_k(\mu_k) \to 0$  in high dimensions (i.e there are few points from Gaussian k near  $\mu_k$ )

(S)

### Other "CS" algorithms I

- Dasgupta (2000) round, equal sized Gaussian, random projection
- Arora and Kannan (2001) arbitrary shaped Gaussian, distances
- Achlioptas and McSherry (2005) log-concave, principal subspace projection

**Example Theorem** (Achlioptas & McSherry, 2005) If data come from K Gaussians,  $n >> K(d + \log K)/\pi_{min}$ , and

$$||\mu_k - \mu_l|| \ge 4\sigma_k \sqrt{1/\pi_k + 1/\pi_l} + 4\sigma_k \sqrt{K\log nK + K^2}$$

then, w.h.p.  $1 - \delta(d, K, n)$ , their algorithm finds true labels **Good** 

- theoretical guarantees
- no local optima
- suggest heuritics for EM K-means
  - project data on principal subspace (when d >> K)

#### But

• strong assuptions: large separation (unrealistic), concentration of  $f_k$ 's (or  $f_k$  known), K known

イロト 不得 トイヨト イヨト

CSE 547/STAT 548 Winter 2022

23

• try to find perfect solution (too ambitious)

### A fundamental result

The Johnson-Lindenstrauss Lemma For any  $\varepsilon \in (0, 1]$  and any integer n, let d' be a positive integer such that  $d' \geq 4(\varepsilon^2/2 - \varepsilon^3/3)^{-1} \ln n$ . Then for any set  $\mathcal{D}$  of n points in  $\mathbb{R}^d$ , there is a map  $f : \mathbb{R}^d \to \mathbb{R}^{d'}$  such that for all  $u, v \in V$ ,

$$(1-\varepsilon)||u-v||^{2} \le ||f(u)-f(v)||^{2} \le (1+\varepsilon)||u-v||^{2}$$
(7)

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

CSE 547/STAT 548 Winter 2022

24

Furthermore, this map can be found in randomized polynomial time.

- note that the embedding dimension d' does not depend on the original dimension d, but depends on  $n, \varepsilon$
- ? show that: the mapping f is linear and that w.p.  $1 \frac{1}{n}$  a random projection (rescaled) has this property
- their proof is elementary Projecting a fixed vector v on a a random subspace is the same as projecting a random vector v on a fixed subspace. Assume  $v = [v_1, \ldots, v_d]$  with  $v \sim i.i.d.$  and let  $\tilde{v} = \text{projection of } v$  on axes 1 : d'. Then  $E[||\tilde{v}||^2 = d'E[v_j^2] = \frac{d'}{d}E[||v||^2]$ . The next step is to show that the variance of  $||\tilde{v}||^2$  is very small when d' is sufficiently large.

# A two-step EM algorithm Dasgupta and Schulman (2007)

Assumes K spherical gaussians, separation  $||\mu_k^{true} - \mu_{k'}^{true} \ge C\sqrt{d\sigma_k}$ 

- Pick  $K' = \mathcal{O}(K \ln K)$  centers  $\mu_k^0$  at random from the data
- (a) Set  $\sigma_k^0 = \frac{d}{2} \min_{k \neq k'} ||\mu_k^0 \mu_{k'}^0||^2$ ,  $\pi_k^0 = 1/K'$
- **(a)** Run one E step and one M step  $\Longrightarrow \{\pi_k^1, \mu_k^1, \sigma_k^1\}_{k=1:K'}$
- Compute "distances"  $d(\mu_k^1, \mu_{k'}^1) = \frac{||\mu_k^1 \mu_{k'}^1||}{\sigma_k^1 \sigma_{k'}^1}$
- **③** Prune all clusters with  $\pi_k^1 \leq 1/4K'$
- Run Fastest First Traversal with distances d(μ<sup>1</sup><sub>k</sub>, μ<sup>1</sup><sub>k'</sub>) to select K of the remaining centers. Set π<sup>1</sup><sub>k</sub> = 1/K.
- $\textbf{0} \ \, {\rm Run \ one \ } {\rm E \ step \ and \ one \ } {\rm M \ step} \Longrightarrow \{\pi_k^2, \mu_k^2, \sigma_k^2\}_{k=1:K}$

eorem For any  $\delta, \varepsilon > 0$  if d large, n large enough, separation  $C \ge d^{1/4}$  the Two step EM algorithm obtains centers  $\mu_k$  so that

$$||\mu_k - \mu_k^{true}|| \le ||\text{mean}(C_k^{true}) - \mu_k^{true}|| + \varepsilon \sigma_k \sqrt{d}$$

イロン イ団 とく ヨン イヨン

CSE 547/STAT 548 Winter 2022

э

26

## Selecting K for mixture models

#### The BIC (Bayesian Information) Criterion

- let  $\theta_K$  = parameters for  $\gamma_K$
- let  $\#\theta_K$ =number independent parameters in  $\theta_K$ 
  - e.g for mixture of Gaussians with full  $\Sigma_k$ 's in d dimensions

$$\#\theta_{K} = \underbrace{K-1}_{\pi_{1:K}} + \underbrace{Kd}_{\mu_{1:K}} + \underbrace{Kd(d-1)/2}_{\Sigma_{1:K}}$$

define

$$BIC(\theta_{K}) = I(\theta_{K}) - \frac{\#\theta_{K}}{2} \ln n$$

- Select K that maximizes  $BIC(\theta_K)$
- selects true K for  $n \to \infty$  and other technical conditions (e.g parameters in compact set)
- but theoretically not justified (and overpenalizing) for finite n

Number of Clusters vs. BIC EII (A), VII (B), EEI (C), VEI (D), EVI (E), VVI (F), EEE (G), EEV (H), VEV (I), VVV (J)





(from Nugent and Meila (2010))

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

Number of Clusters vs. BIC EII (A), VII (B), EEI (C), VEI (D), EVI (E), VVI (F), EEE (G), EEV (H), VEV (I), VVV (J)

EEV, 8 Cluster Solution

イロト イヨト イヨト イヨ

CSE 547/STAT 548 Winter 2022

28



(from Nugent and Meila (2010))

Marina Meila (UW)

### Selecting K for hard clusterings

- based on statistical testing: the gap statistic (Tibshirani, Walther, Hastie, 2000)
- the Krzanowski-Lai (KL), 1985 statistic
- X-means Pelleg and Moore (2000) heuristic: splits/merges clusters based on statistical tests of Gaussianity

- Stability methods
  - Empirical prove instability
  - Optimization based prove stability

### Empirical Stability methods for choosing K

- like bootstrap, or crossvalidation
- Idea (implemented by Ben-Hur et al. (2002))
  - for each K

    - (2) cluster  $\mathcal{D}' \to \Delta'_K$
    - compare Δ<sub>K</sub>, Δ'<sub>K</sub>. Are they similar?
       If yes, we say Δ<sub>K</sub> is stable to perturbations

Fundamental assumption If  $\Delta_K$  is stable to perturbations then K is the correct number of clusters

- these methods are supported by experiments (not extensive)
- not directly supported by theory ... see von Luxburg (2009) for a summary of the area

#### The gap statistic

#### Idea

- for some cost L compare  $L(\Delta_K)$  with its expected value under a null distribution
  - · choose null distribution to have no clusters
    - Gaussian (fit to data)
    - uniform with convex support
    - uniform over K<sub>0</sub> principal components of data
  - null value =  $E_{P_0}[L_{K,n}]$  the expected value of the cost of clustering *n* points from  $P_0$  into *K* clusters
- the gap

$$g(K) = E_{P_0}[L_{K,n}] - L(\Delta_K) = L_K^0 - L(\Delta_K)$$

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

3

- choose  $K^*$  corresponding to the largest gap
- nice: it can also indicate that data has no clusters

#### Practicalities

- $L_{K}^{0} = E_{P_{0}}[L_{K,n}]$  can rarely be computed in closed form (when  $P_{0}$  very simple)
- otherwise, estimate L<sup>0</sup><sub>K</sub> be Monte-Carlo sampling i.e generate B samples from P<sub>0</sub> and cluster them
- if sampling, variance  $s_K^2$  of estimate  $\hat{L}_K^0$  must be considered  $s_K^2$  is also estimated from the samples
- selection rule:  $K^* = \text{smallest } K$  such that  $g(K) \geq g(K+1) s_{K+1}$
- favored  $L^V(\Delta) = \sum_k \frac{1}{|C_k|} \sum_{i \in C_k} ||x_i \mu_k||^2 \approx \text{sum of cluster variances}$

イロト 不得 トイヨト イヨト 二日

#### The KL statistic

- Heuristic
- define  $w_{\mathcal{K}} = \mathcal{K}^{2/d} L^{\mathcal{V}}(\Delta_{\mathcal{K}})$  (lower is better)
- a good K is indicated by "large" jump between K 1 and K
- they propose

$$diff(K) = w_{K-1} - w_K$$

choose K that maximixes relative jump  $\left| \frac{diff(K)}{diff(K+1)} \right|$ 

34

## Stability methods for choosing K

- like bootstrap, or crossvalidation
- Idea (implemented by Ben-Hur et al. (2002))

for each K

- $\ \, \bullet \ \, \bullet \ \, \bullet \ \, \bullet \ \, {\mathcal D}'$
- (2) cluster  $\mathcal{D}' \to \Delta'_K$
- Socompare  $\Delta_K$ ,  $\Delta'_K$ . Are they similar? If yes, we say  $\Delta_K$  is stable to perturbations

Fundamental assumption If  $\Delta_K$  is stable to perturbations then K is the correct number of clusters

- these methods are supported by experiments (not extensive)
- not YET supported by theory ... see von Luxburg (2009) for a summary of the area

イロト イヨト イヨト イヨト

CSE 547/STAT 548 Winter 2022

э

35

## A stability based method for model-based clustering

#### • The algorithm of Lange et al. (2004)

- ( )divide data into 2 halves  $\mathcal{D}_1, \mathcal{D}_2$  at random
- 2 cluster (by EM)  $\mathcal{D}_1 \rightarrow \Delta_1, \theta_1$
- **③** cluster (by EM)  $D_2$  →  $Δ_2$ ,  $θ_2$
- **(**) compare  $\Delta_1, \Delta'_1$
- o repeat B times and average the results
- repeat for each  ${\cal K}$
- select K where Δ<sub>1</sub>, Δ'<sub>1</sub> are closest on average (or most times)



Fig. 2.1 Normalized stability scores. Left plots: Data points from a uniform density on  $[0,1]^2$ . Right plots: Data points from a mixture of four well-separated Gaussians in  $\mathbb{R}^2$ . The first row always shows the unnormalized instability Instab for K = 2, ..., 15. The second row shows the instability Instab<sub>norm</sub> obtained on a reference distribution (uniform distribution). The third row shows the normalized stability Instab<sub>norm</sub>.

イロト イヨト イヨト

CSE 547/STAT 548 Winter 2022

36

(from von Luxburg (2009))

### Clustering with outliers

- What are outliers?
- let p = proportion of outliers (e.g 5%-10%)
- Remedies
  - mixture model: introduce a K + 1-th cluster with large (fixed)  $\Sigma_{K+1}$ , bound  $\Sigma_k$  away from 0

- K-means and EM
  - robust means and variances
    - e.g eliminate smallest and largest  $pn_k/2$  samples in mean computation (trimmed mean)
  - K-medians Charikar and Guha (1999)
  - replace Gaussian with a heavier-tailed distribution (e.g. Laplace)
- single-linkage: do not count clusters with < r points
- Is K meaningful when outliers present?
  - alternative: non-parametric clustering

#### References

- Achlioptas, D. and McSherry, F. (2005). On spectral learning of mixtures of distributions. In Auer, P. and Meir, R., editors, 18th Annual Conference on Learning Theory, COLT 2005, pages 458–471, Berlin/Heidelberg. Springer.
- Arora, S. and Kannan, R. (2001). Learning mixtures of arbitrary gaussians. In STOC '01: Proceedings of the thirty-third annual ACM symposium on Theory of computing, pages 247–257, New York, NY, USA. ACM Press.
- Banfield, J. D. and Raftery, A. E. (1993). Model-based gaussian and non-gaussian clustering. Biometrics, 49:803–821.
- Ben-Hur, A., Elisseeff, A., and Guyon, I. (2002). A stability based method for discovering structure in clustered data. In *Pacific Symposium on Biocomputing*, pages 6–17.
- Carreira-Perpinan, M. A. (2007). Gaussian mean shift is an EM algorithm. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 29(5):767–776.
- Charikar, M. and Guha, S. (1999). Improved combinatorial algorithms for the facility location and k-median problems. In 40th Annual Symposium on Foundations of Computer Science, pages 378–388.
- Dasgupta, S. (2000). Experiments with random projection. In UAI '00: Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence, pages 143–151, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Dasgupta, S. and Gupta, A. (2002). An elementary proof of a theorem of johnson and lindenstrauss. *Algorithms*, 22:60–65.
- Dasgupta, S. and Schulman, L. (2007). A probabilistic analysis of em for mixtures of separated, spherical gaussians. *Journal of Machine Learnig Research*, 8:203–226.

Marina Meila (UW)

CSE 547/STAT 548 Winter 2022 37

- Lange, T., Roth, V., Braun, M. L., and Buhmann, J. M. (2004). Stability-based validation of clustering solutions. Neural Comput., 16(6):1299-1323.
- McLachlan, G. J. and Krishnan, T. (1997). The EM algorithm and extensions. Wiley, New York, NY.
- Neal, R. M. and Hinton, G. E. (1998). A view of the em algorithm that justifies incremental, sparse, and other variants. In Jordan, M. I., editor, Learning in Graphical Models, NATO Science series, pages 355-368. Kluwer Academic Publishers.
- Nugent, R. and Meila, M. (2010). Statistical Methods in Molecular Biology, chapter An Overview of Clustering Applied to Molecular Biology. Humana Press, Springer.
- Pelleg, D. and Moore, A. (2000). X-means: Extending K-means with efficient estimation of the number of clusters. In Bratko, I. and Džeroski, S., editors, Proceedings of the 17th International Conference on Machine Learning, pages 727–734, San Francisco, CA. Morgan Kaufmann.
- Vempala, S. and Wang, G. (2004). A spectral algorithm for learning mixtures of distributions. Journal of Computer Systems Science, 68(4):841–860.
- von Luxburg, U. (2009). Clustering stability: An overview. Foundations and Trends in Machine Learning, 2(3):253-274.

イロト 不得 トイヨト イヨト 二日

CSE 547/STAT 548 Winter 2022