

Advanced Clustering Techniques

Cluster Analysis: Advanced Methods

- Probability Model-Based Clustering
- Clustering High-Dimensional Data
- Clustering with Constraints
- Summary



Need for advanced cluster Analysis

- *Online store : customer purchases, create reviews of products*
- *Review may involve multiple products and if we want to cluster the review*
- *Assigning a review to a cluster exclusively would not work for our task*
- *We need a clustering method that allows a review that belong to more than one cluster if the review indeed involves more than a topic.*
- *Assignment of a review to any cluster involve determination of weight representing the partial membership*

Need for advanced cluster Analysis

- *Clustering methods discussed so far*
 - *Every data object is assigned to exactly one cluster*
- *Some applications may need for fuzzy or soft cluster assignment*
 - *Ex. An e-game could belong to both entertainment and software*
 - *Methods: fuzzy clusters and probabilistic model-based clusters*

Fuzzy Set and Fuzzy Cluster

- Given a set of objects, $X = \{x_1, \dots, x_n\}$, a fuzzy set S is a subset of X that allows each object in X to have a membership degree between 0 and 1.
- Fuzzy cluster: A fuzzy set $S: F_S: X \rightarrow [0, 1]$ (value between 0 and 1)
- Example: Popularity of cameras is defined as a fuzzy mapping
- Function $\text{pop}()$ defines a fuzzy set for digital cameras
- $\{A(0.05), B(1), C(0.86), D(0.27)\}$ where the degrees of membership

Camera	Sales (units)
A	50
B	1320
C	860
D	270

$$\text{Pop}(o) = \begin{cases} 1 & \text{if 1,000 or more units of } o \text{ are sold} \\ \frac{i}{1000} & \text{if } i \text{ } (i < 1000) \text{ units of } o \text{ are sold} \end{cases}$$

Fuzzy Set and Fuzzy Cluster

- Given a set of objects, a cluster is a fuzzy set of objects. Such a cluster is called a fuzzy cluster
- Given a set of objects, o_1, \dots, o_n , a fuzzy clustering of k fuzzy clusters, C_1, \dots, C_k , can be represented using a partition matrix, $M = [W_{ij}]$ ($1 \leq i \leq n, 1 \leq j \leq k$) where W_{ij} is the membership of degree of o_i in cluster C_j .
- The matrix should follow three requirements :
 - P1: for each object o_i and cluster C_j , $0 \leq w_{ij} \leq 1$ (fuzzy set)
 - P2: for each object o_i , $\sum_{j=1}^k w_{ij} = 1$, equal participation in the clustering
 - P3: for each cluster C_j , $0 < \sum_{i=1}^n w_{ij} < n$ ensures there is no empty cluster

Fuzzy (Soft) Clustering

■ Example: Let cluster features be

- C_1 : “digital camera” and “lens”
- C_2 : “computer”

Review-id	Keywords
R_1	digital camera, lens
R_2	digital camera
R_3	lens
R_4	digital camera, lens, computer
R_5	computer, CPU
R_6	computer, computer game

$$M = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ \frac{2}{3} & \frac{1}{3} \\ 0 & 1 \\ 0 & 1 \end{bmatrix}$$

- $W_{ij} = \frac{|R_i \cap C_j|}{|R_i \cap (C_1 \cup C_2)|}$
- Let c_1, \dots, c_k as the center of the k clusters
- For an object o_i , and cluster C_j if $W_{ij} > 0$ then $\text{dist}(o_i, C_j)$ measures how well o_i is represented by c_j , and thus belongs to cluster C_j .
- An object can participate in more than one cluster the sum of distances to the corresponding cluster centers weighted by the degrees of membership captures how well the objects fits the clustering

Fuzzy (Soft) Clustering

- For an object o_i sum of the squared error (SSE), is given by

$$SSE(o_i) = \sum_{j=1}^k w_{ij}^p \text{dist}(o_i, c_j)^2$$

- where p is a parameter $p(p \geq 1)$ controls the influence of the degrees of membership.

- For a cluster C_j , SSE:

$$SSE(C_j) = \sum_{i=1}^n w_{ij}^p \text{dist}(o_i, c_j)^2$$

- Measure how well a clustering fits the

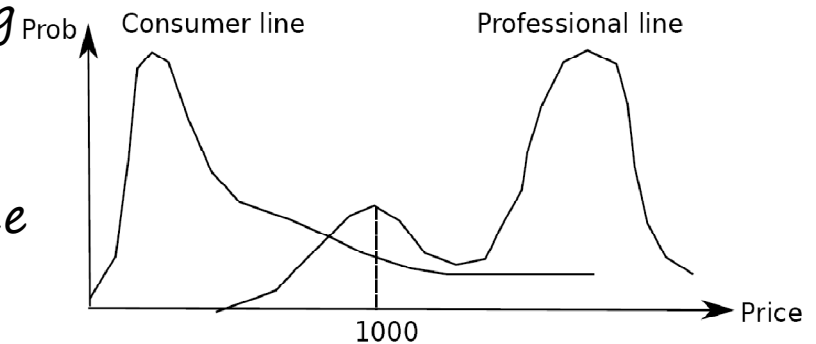
$$SSE(C) = \sum_{i=1}^n \sum_{j=1}^k w_{ij}^p \text{dist}(o_i, c_j)^2$$

Probabilistic Model-Based Clustering

- The inherent categories hidden in the data are latent they cannot be directly observed but we need to infer them.
- The goal of cluster analysis is to find hidden categories.
- A hidden category (i.e., probabilistic cluster) is a distribution over the data space, that are inferred using the data set and are designed to approach hidden categories.
- It can be mathematically represented using a probability density function (or distribution function).
- For a probabilistic cluster, C , its probability density function, f , and a point, o , in the data space, $f(o)$ is the relative likelihood that an instance of C appears at o .

Probabilistic Model-Based Clustering

- Ex. 2 categories for digital cameras sold
 - consumer line vs. professional line
 - density functions f_1, f_2 for C_1, C_2
 - obtained by probabilistic clustering
 - F_1 and f_2 are not observed directly and can be inferred by analyzing the price



- A mixture model assumes that a set of observed objects is a mixture of instances from multiple probabilistic clusters, and conceptually each observed object is generated independently
- Out task: infer a set of k probabilistic clusters that is mostly likely to generate D using the above data generation process

Model-Based Clustering

- A set \mathcal{C} of k probabilistic clusters $\mathcal{C}_1, \dots, \mathcal{C}_k$ with probability density functions f_1, \dots, f_k , respectively, and their probabilities $\omega_1, \dots, \omega_k$.
- Probability of an object o generated by cluster \mathcal{C}_j is $P(o|\mathcal{C}_j) = \omega_j f_j(o)$
- Probability of o generated by the set of cluster \mathcal{C} is $P(o|\mathcal{C}) = \sum_{j=1}^k \omega_j f_j(o)$
- Since objects are assumed to be generated independently, for a data set $D = \{o_1, \dots, o_n\}$, we have,
$$P(D|\mathcal{C}) = \prod_{i=1}^n P(o_i|\mathcal{C}) = \prod_{i=1}^n \sum_{j=1}^k \omega_j f_j(o_i)$$
- Task: Find a set \mathcal{C} of k probabilistic clusters s.t. $P(D|\mathcal{C})$ is maximized
- However, maximizing $P(D|\mathcal{C})$ is often intractable since the probability density function of a cluster can take an arbitrarily complicated form
- To make it computationally feasible (as a compromise), assume the probability density functions being some parameterized distributions

Univariate Gaussian Mixture Model

- $O = \{o_1, \dots, o_n\}$ (n observed objects), $\Theta = \{\theta_1, \dots, \theta_k\}$ (parameters of the k distributions), and $P_j(o_i | \theta_j)$ is the probability that o_i is generated from the j -th distribution using parameter θ_j , we have

$$P(o_i | \Theta) = \sum_{j=1}^k \omega_j P_j(o_i | \theta_j) \quad P(O | \Theta) = \prod_{i=1}^n \sum_{j=1}^k \omega_j P_j(o_i | \theta_j)$$

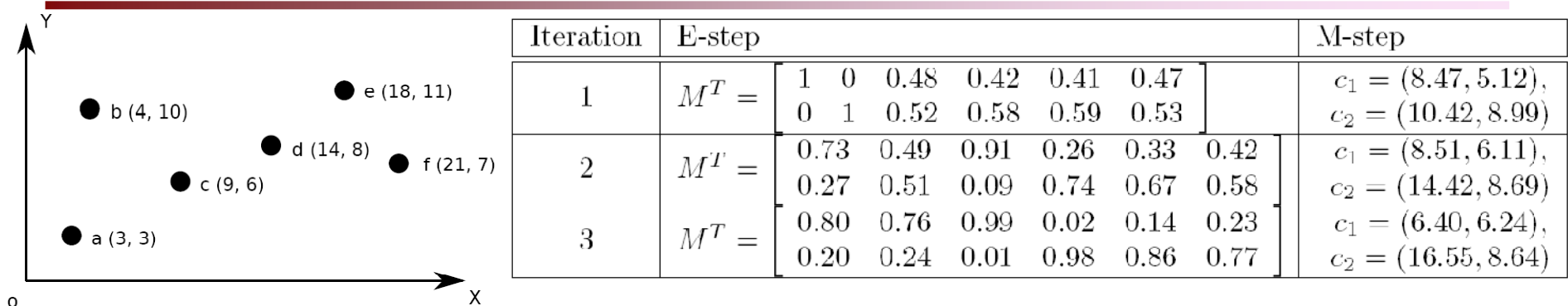
- Univariate Gaussian mixture model
 - Assume the probability density function of each cluster follows a 1-d Gaussian distribution. Suppose that there are k clusters.
 - The probability density function of each cluster are centered at μ_j with standard deviation σ_j , $\theta_j = (\mu_j, \sigma_j)$, we have

$$P(o_i | \theta_j) = \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}} \quad P(o_i | \Theta) = \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$
$$P(O | \Theta) = \prod_{i=1}^n \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$

The EM (Expectation Maximization) Algorithm

- The k-means algorithm has two steps at each iteration:
 - **Expectation Step (E-step):** Given the current cluster centers, each object is assigned to the cluster whose center is closest to the object: An object is *expected to belong to the closest cluster*
 - **Maximization Step (M-step):** Given the cluster assignment, for each cluster, the algorithm *adjusts the center so that the sum of distance from the objects assigned to this cluster and the new center is minimized*
- **The (EM) algorithm:** A framework to approach maximum likelihood or maximum a posteriori estimates of parameters in statistical models.
 - **E-step** assigns objects to clusters according to the current fuzzy clustering or parameters of probabilistic clusters
 - **M-step** finds the new clustering or parameters that maximize the sum of squared error (SSE) or the expected likelihood

Fuzzy Clustering Using the EM Algorithm



- Initially, let $c_1 = a$ and $c_2 = b$

- 1st E-step: assign o to c_1 , w. wt = $\frac{\frac{1}{dist(o, c_1)^2}}{\frac{1}{dist(o, c_1)^2} + \frac{1}{dist(o, c_2)^2}} = \frac{dist(o, c_2)^2}{dist(o, c_1)^2 + dist(o, c_2)^2}$

$$w_{c, c_1} = \frac{41}{45+41} = 0.48$$

- 1st M-step: recalculate the centroids according to the partition matrix, minimizing the sum of squared error (SSE)

$$c_j = \frac{\sum_{\text{each point } o} w_{o, c_j}^2 o}{\sum_{\text{each point } o} w_{o, c_j}^2} \quad c_1 = \left(\frac{1^2 \times 3 + 0^2 \times 4 + 0.48^2 \times 9 + 0.42^2 \times 14 + 0.41^2 \times 18 + 0.47^2 \times 21}{1^2 + 0^2 + 0.48^2 + 0.42^2 + 0.41^2 + 0.47^2}, \frac{1^2 \times 3 + 0^2 \times 10 + 0.48^2 \times 6 + 0.42^2 \times 8 + 0.41^2 \times 11 + 0.47^2 \times 7}{1^2 + 0^2 + 0.48^2 + 0.42^2 + 0.41^2 + 0.47^2} \right) = (8.47, 5.12)$$

- Iteratively calculate this until the cluster centers converge or the change is small enough

Univariate Gaussian Mixture Model

- $O = \{o_1, \dots, o_n\}$ (n observed objects), $\Theta = \{\theta_1, \dots, \theta_k\}$ (parameters of the k distributions), and $P_j(o_i | \theta_j)$ is the probability that o_i is generated from the j-th distribution using parameter θ_j , we have

$$P(o_i | \Theta) = \sum_{j=1}^k \omega_j P_j(o_i | \theta_j) \quad P(O | \Theta) = \prod_{i=1}^n \sum_{j=1}^k \omega_j P_j(o_i | \theta_j)$$

- Univariate Gaussian mixture model
 - Assume the probability density function of each cluster follows a 1-d Gaussian distribution. Suppose that there are k clusters.
 - The probability density function of each cluster are centered at μ_j with standard deviation σ_j , $\theta_j = (\mu_j, \sigma_j)$, we have

$$P(o_i | \theta_j) = \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}} \quad P(o_i | \Theta) = \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$
$$P(O | \Theta) = \prod_{i=1}^n \sum_{j=1}^k \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(o_i - \mu_j)^2}{2\sigma_j^2}}$$

Computing Mixture Models with EM

- Given n objects $\mathbf{O} = \{o_1, \dots, o_n\}$, we want to mine a set of parameters $\Theta = \{\theta_1, \dots, \theta_k\}$ s.t., $P(\mathbf{O}|\Theta)$ is maximized, where $\theta_j = (\mu_j, \sigma_j)$ are the mean and standard deviation of the j -th univariate Gaussian distribution
- We initially assign random values to parameters θ_j , then iteratively conduct the E- and M- steps until converge or sufficiently small change
- At the E-step, for each object o_i , calculate the probability that o_i belongs to each distribution,

$$P(\theta_j|o_i, \Theta) = \frac{P(o_i|\theta_j)}{\sum_{l=1}^k P(o_i|\theta_l)}$$

- At the M-step, adjust the parameters $\theta_j = (\mu_j, \sigma_j)$ so that the expected likelihood $P(\mathbf{O}|\Theta)$ is maximized

$$\mu_j = \sum_{i=1}^n o_i \frac{P(\theta_j|o_i, \Theta)}{\sum_{l=1}^k P(\theta_j|o_l, \Theta)} = \frac{\sum_{i=1}^n o_i P(\theta_j|o_i, \Theta)}{\sum_{i=1}^n P(\theta_j|o_i, \Theta)} \quad \sigma_j = \sqrt{\frac{\sum_{i=1}^n P(\theta_j|o_i, \Theta)(o_i - \mu_j)^2}{\sum_{i=1}^n P(\theta_j|o_i, \Theta)}}$$

Advantages and Disadvantages of Mixture Models

- Strength
 - Mixture models are more general than partitioning and fuzzy clustering
 - Clusters can be characterized by a small number of parameters
 - The results may satisfy the statistical assumptions of the generative models
- Weakness
 - Converge to local optimal (overcome: run multi-times w. random initialization)
 - Computationally expensive if the number of distributions is large, or the data set contains very few observed data points
 - Need large data sets
 - Hard to estimate the number of clusters