

CLUSTER ANALYSIS

1 Introduction

Reference: Chapter 12, Johnson and Wichern

Clustering methods were developed to detect inherent structure of the data.

There number of clustering methods available to group multivariate data.

Assumption :

The only assumption required is number of inherent clusters in the data.

That is, assumption on distribution is not required to apply any clustering method. Therefore, it is a distribution free (non-parametric) procedure.

2 Minimum Distance Clustering

Let $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p})^T (i = 1, 2, \dots, n)$ and $\mathbf{x}_j = (x_{j,1}, x_{j,2}, \dots, x_{j,p})^T (j = 1, 2, \dots, n)$ be two p -dimensional observations.

Definitions:

(a) Euclidean Distance

$$\begin{aligned} D_{i,j} &= \|\mathbf{x}_i - \mathbf{x}_j\| = \left\{ (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j) \right\}^{1/2} \\ &= \left\{ \sum_{s=1}^p (x_{i,s} - x_{j,s})^2 \right\}^{1/2} \end{aligned}$$

(b) L_l Distance

$$L_{i,j} = |\mathbf{x}_i - \mathbf{x}_j| = \sum_{s=1}^p |x_{i,s} - x_{j,s}|$$

Now group the observations so that the total distance (Euclidean or L_1) within the group is minimized.

L_1 distance is computationally faster to determine but it is less accurate than the Euclidean distance method.

3 Clustering Quality Indicator

A common clustering criterion or quality indicator is the sum of squares of errors

$$\begin{aligned}\text{SSE} &= \sum_{r=1}^k \sum_{\mathbf{x}_t \in C_r} \|\mathbf{x}_t - \boldsymbol{\mu}_r\|^2 \\ &= \sum_{r=1}^k \sum_{\mathbf{x}_t \in C_r} (\mathbf{x}_t - \boldsymbol{\mu}_r)^T (\mathbf{x}_t - \boldsymbol{\mu}_r) \\ &= \sum_{r=1}^k \sum_{\mathbf{x}_t \in C_r} \sum_{s=1}^p (x_{t,s} - \mu_{r,s})^2\end{aligned}$$

Where $\boldsymbol{\mu}_r = (\mu_{r,1}, \mu_{r,2}, \dots, \mu_{r,p})$ is the mean of cluster C_r ($r = 1, \dots, k$) and $\mathbf{x}_t \in C_r$ implies that observation \mathbf{x}_t is assigned to cluster C_r .

4 Migrating Means Algorithm (K Mean Method)

Step 1: Assume there are k clusters in the multivariate space. Select k different points arbitrary, say, $\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_k$ from the space to initialize the k means of the clusters. Or divide \mathbf{x}_i 's into k clusters arbitrary and compute the average of \mathbf{x}_i 's in C_r and let \mathbf{m}_r be the average of \mathbf{x}_i 's in C_r .

Step 2: Compute EITHER Euclidean distance, $D_{i,j}$ or L_1 distance where

$$\begin{aligned}D_{i,j} &= \|\mathbf{x}_i - \mathbf{m}_j\| = \left\{ (\mathbf{x}_i - \mathbf{m}_j)^T (\mathbf{x}_i - \mathbf{m}_j) \right\}^{1/2} \\ &= \left\{ \sum_{s=1}^p (x_{i,s} - m_{j,s})^2 \right\}^{1/2}\end{aligned}$$

and

$$L_{ij} = |\mathbf{x}_i - \mathbf{m}_j| = \sum_{s=1}^p |x_{i,s} - m_{j,s}|$$

where $\mathbf{m}_j = (m_{j,1}, m_{j,2}, \dots, m_{j,p})^T$.

Step 3: Assign the observation, \mathbf{x}_i to the nearest cluster for all i .

That is, assign \mathbf{x}_i to C_r if $D_{i,r} = \min_j (D_{i,j})$ or for L_1 distance measure $L_{i,r} = \min_j (L_{i,j})$.

Step 4: Compute the new mean for clusters using the assigned observations. Let $\hat{\boldsymbol{\mu}}_r = \frac{1}{n_r} \sum_{\mathbf{x}_t \in C_r} \mathbf{x}_t$ where n_r is the number of observations assigned to C_r .

If $\hat{\boldsymbol{\mu}}_r = \mathbf{m}_r$ for all r then the assignment of observations to clusters are completed and **STOP** the procedure. Otherwise, set $\mathbf{m}_r = \hat{\boldsymbol{\mu}}_r$ and GOTO Step 2.

Example: Assign the following 4 observations into two groups:

$$\begin{pmatrix} 5 \\ 3 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -2 \end{pmatrix} \text{ and } \begin{pmatrix} -3 \\ -2 \end{pmatrix}.$$

5 Hierarchical Clustering

This method does not require the user to specify the number of clusters. It produces an output that allows the user to decide the natural clusters in the population.

Initially, hierarchical clustering method assumes that all objects are individual clusters.

Next merges the neighbouring pair of clusters by checking the Distance (Similarity) matrix. Continue this until all objects appear in a single cluster.

Procedure

Step 1: Let N be the number of objects in the multivariate population and $\mathcal{D} = (d_{i,j})_{N \times N}$ be the distance matrix (dissimilarity matrix), where $d_{i,j}$ is the distance (dissimilarity) between the objects (or observations) i and j .

Note that $d_{i,j}$ can be an Euclidean distance, L_1 distance or $1 - r_{i,j}^2$, where $r_{i,j}$ is the sample correlation between i and j , or any other similar quantity.

Step 2: Search the distance matrix for the most similar (nearest) pair of clusters. Let u and v be the most similar clusters(objects) and uv be the new cluster formed by merging u and v .

Step 3: Delete the rows and columns of the distance matrix corresponding to clusters u and v and add a row and column giving the distances between cluster uv and the remaining clusters.

Repeat the Steps 2 and 3 a total of $N - 1$ times. Then all objects will be in a single cluster. Record the identity of clusters that are merged and the levels and present the results graphically in the form of a **dendrogram**.

Linkage Methods

That is, methods of combining two or more objects to form a cluster. Commonly used linkage method are:

- Single Linkage
- Complete Linkage
- Average Linkage

Single Linkage (Nearest Neighbour) Method

Find the smallest distance in \mathcal{D} and merge the corresponding objects, say u and v , to get the cluster uv . For the above Step 3, the distances between uv and another cluster(or object) w (say), is computed by

$$d_{uv,w} = \min(d_{u,w}, d_{v,w})$$

The other linkage methods proceeds in much the same way as single linkage, with one exception. That is, the computation of the distance between newly form cluster and other clusters(or objects).

Note that

- For **complete Linkage (Furthest neighbour) Method:**

$$d_{uv,w} = \max(d_{u,w}, d_{v,w})$$

- For **Average Linkage Method:**

$$d_{uv,w} = \frac{1}{N_{uv}N_w} \sum_i \sum_j d_{i,j}$$

where $d_{i,j}$ is the distance between object i in cluster uv and object j in cluster w and, N_{uv} and N_w are the number of objects in clusters uv and w respectively.

Example: The distances between pairs of five item are given below:

	1	2	3	4	5
1	0				
2	4	0			
3	6	9	0		
4	1	7	10	0	
5	6	3	5	8	0

Cluster the five items using the

- Single linkage
- Complete linkage and
- Average linkage

methods. Draw a dendograms.

Distance Measures for Binary Variables

Let $x_{i,k}$ be the score (1 or 0) of binary variable X_k on i^{th} object \mathbf{x}_i and $x_{j,k}$ be the score (1 or 0) of X_k on j^{th} object for $k = 1, 2, \dots, p$.

Let

- $a = (1 - 1)$ matches between object i and j ,
- $b = (1 - 0)$ mismatches between object i and j ,
- $c = (0 - 1)$ mismatches between object i and j , and
- $d = (0 - 0)$ matches between object i and j .

Then, $a + b + c + d = p$.

Example: Certain characteristics associated with six workers in a factory are listed below:

Person	Married	School	Sex
1	yes	Private	Female
2	no	state	Female
3	yes	state	Male
4	yes	state	Male
5	no	Private	Female
6	no	Private	Male

Draw a dendograms and identify similar groups (clusters) among these persons.

The following tables (Table 12.1, Johnson and Wichern, p.675) lists commonly use similarity coefficients for bivariate data.

Table: Similarity Coefficients For Clustering

Coefficient	Rationale
1. $\frac{a+d}{p}$	Equal weights for 1-1 matches and 0-0 matches.
2. $\frac{2(a+d)}{2(a+d)+b+c}$	Double weight for 1-1 matches and 0-0 matches.
3. $\frac{a+d}{a+d+2(b+c)}$	Double weight for unmatched pairs.
4. $\frac{a}{p}$	No 0-0 matches in numerator.
5. $\frac{a}{a+b+c}$	No 0-0 matches in numerator or denominator. The 0-0 matches are treated as irrelevant.
6. $\frac{2a}{2a+b+c}$	No 0-0 matches in numerator or denominator. Double weight for 1-1 matches.
7. $\frac{a}{a+2(b+c)}$	No 0-0 matches in numerator or denominator. Double weight for unmatched pairs.
8. $\frac{a}{b+c}$	Ratio of matches to mismatches with 0-0 matches excluded.

Note that the similarity coefficients 1, 2 and 3 are monotonic.