

به نام خدا

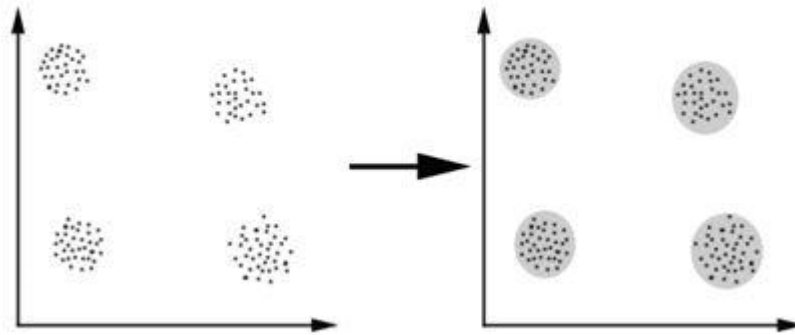
خوشه بندی

CLUSTERING

پاییز ۹۶

مقدمه:

داده و الگو یکی از شاخص‌های بسیار مهم در دنیای اطلاعات هستند. خوشه بندی یکی از بهترین روش‌هایی است که برای کار با داده‌ها ارائه شده است. قابلیت آن در ورود به فضای داده و تشخیص ساختار آنها، خوشه بندی را یکی از ایده‌آل‌ترین مکانیزم‌ها برای کار با دنیای عظیم داده‌ها کرده است. اولین بار ایده‌ی آن در سال ۱۹۳۵ ارائه شد و امروزه با پیشرفت‌ها و جهش‌های عظیمی که در آن پدید آمده، خوشه بندی در کاربردها و جنبه‌های مختلفی حضور یافته است. خوشه بندی یکی از شاخه‌های یادگیری بدون نظارت می‌باشد و فرایند خودکاری است که در طی آن، نمونه‌ها به دسته‌هایی که اعضای آن مشابه یکدیگر می‌باشند تقسیم می‌شوند که به این دسته‌ها خوشه گفته می‌شود. بنابراین **خوشه مجموعه‌ای از اشیاء می‌باشد که در آن اشیاء با یکدیگر مشابه بوده و با اشیاء موجود در خوشه‌های دیگر غیر مشابه می‌باشند.** برای مشابه بودن می‌توان معیارهای مختلفی را در نظر گرفت مثلاً می‌توان معیار فاصله را برای خوشه بندی مورد استفاده قرار داد و اشیائی را که به یکدیگر نزدیکتر هستند را بعنوان یک خوشه در نظر گرفت که به این نوع خوشه بندی، خوشه بندی مبتنی بر فاصله نیز گفته می‌شود. بعنوان مثال در شکل ۱ نمونه‌های ورودی در سمت چپ به چهار خوشه مشابه شکل راست تقسیم می‌شوند. در این مثال هر یک از نمونه‌های ورودی به یکی از خوشه‌ها تعلق دارد و نمونه‌ای وجود ندارد که متعلق به بیش از یک خوشه باشد.



شکل ۱ در این شکل نمونه‌ای از اعمال خوشه‌بندی روی یک مجموعه از داده‌ها مشخص شده است که از معیار فاصله به عنوان عدم شباهت بین داده‌ها استفاده شده است.

What is Clustering?

- Find **K** clusters (or a classification that consists of **K** clusters) so that the objects of one cluster are similar to each other whereas objects of different clusters are dissimilar.
- **Clustering** is a **process of partitioning** a set of data (or objects) in a set of meaningful sub-classes, called **clusters**.

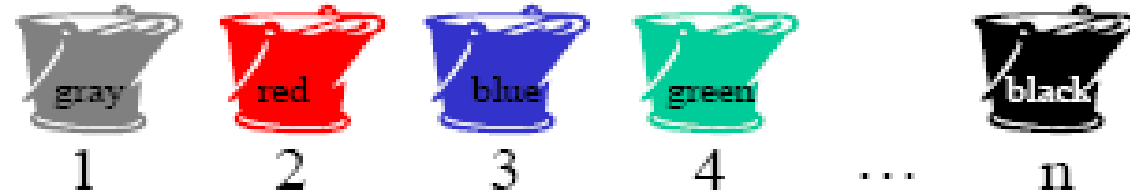
Clustering: unsupervised classification: no predefined classes.

هر کلاستر (خوشه) توسط خصوصیات مشترک اشیائی که درون آن قرار می گیرند تعریف می شود.
کلاستر، یک مجموعه از داده هاست بطوریکه:
۱- داده های موجود در یک کلاستر شبیه یکدیگر هستند.
۲- داده های موجود در کلاسترهای مختلف به یکدیگر شبیه نیستند.

Supervised and Unsupervised

Supervised Classification = Classification

→ We know the class labels and the number of classes



Unsupervised Classification = Clustering

→ We do not know the class labels and may not know the number of classes



Types of Features

- continuous range
- Finite discrete set.
- *Nominal*
 - sex of an individual
 - 1 for a male and 0 for a female.
- *Ordinal*
 - Values can be *meaningfully ordered*
 - Ex. performance of a student in a class
 - are 4,3,2,1 = “excellent,” “very good,” “good,” “not good.”

What is a good clustering?

Internal (within the cluster) distances should be small

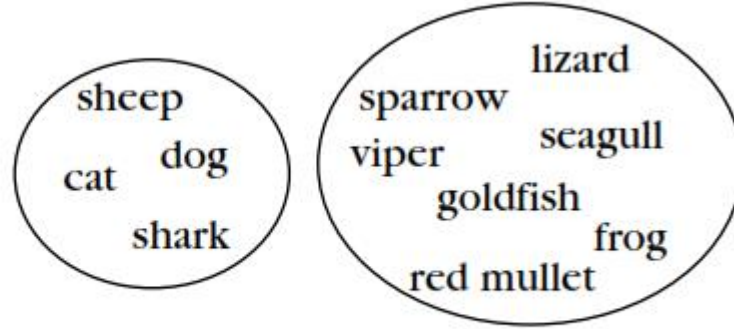
External (intra-cluster) should be large

Clustering is a way to discover new categories (classes)

محاسبه فاصله بین دو داده در خوشه بندی بسیار مهم هست. فاصله که همان معرف عدم تجانس است به ما کمک می کند در فضای داده ای حرکت کنیم و خوشه ها را تشکیل دهیم. با محاسبه فاصله بین دو داده می توان فهمید که چقدر این دو داده به هم نزدیک هستند و بر این اساس آن ها را در یک خوشه قرار داد. توابع ریاضی مختلفی برای محاسبه فاصله وجود دارند. مانند فاصله اقلیدسی، همینگ و...

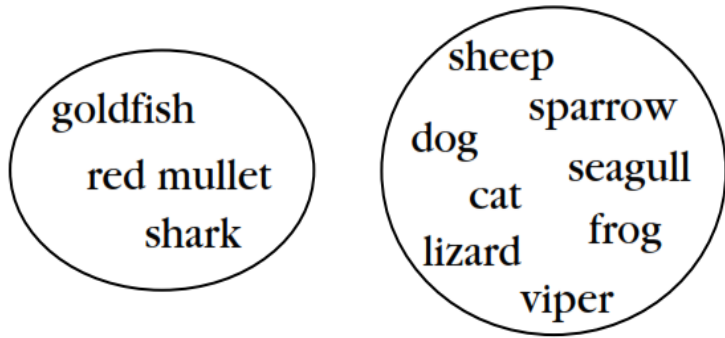
مثال زیر از زیست شناسی بر گرفته شده است و صورت مسئله را برای ما واضح می سازد. به حیوانات زیر توجه کنید: گوسفند، سگ و گربه (پستاندار)، گنجشک و مرغ دریایی (پرنده)، ماهی قرمز، شاه ماهی (ماهی)، افعی و مارمولک (خزنده) و قورباغه (دوزیست). به منظور مرتب کردن این حیوانات در داخل خوشه ها نیاز داریم که یک معیار دسته بندی تعریف کنیم.

اگر ملاک دسته بندی نحوه بدنیا آوردن فرزندان باشد کلاستر (a).



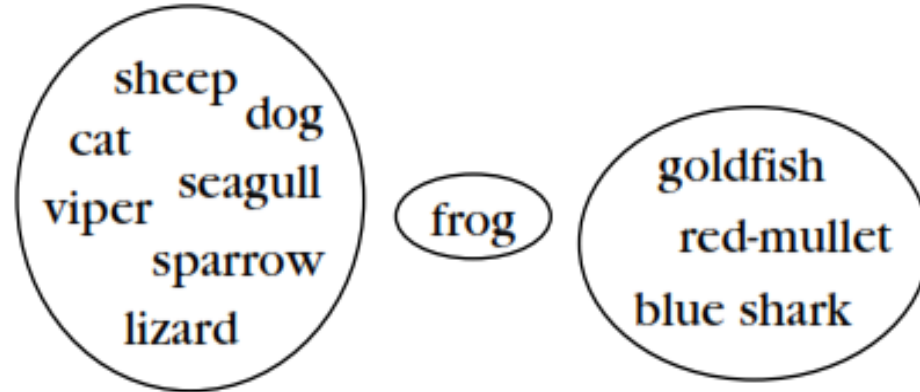
(a)

اگر وجود شش ها را بررسی کنیم، ماهی قرمز، شاه ماهی و کوسه در یک کلاستر و بقیه در یک کلاستر دیگر قرار می گیرند. (b)



(b)

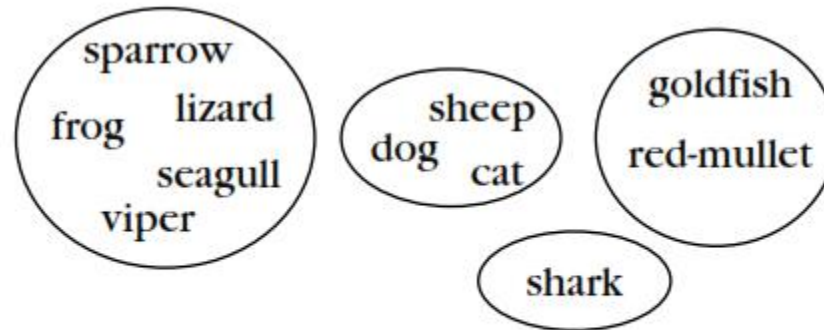
اگر ملاک دسته بندی را محیطی که حیوانات زندگی می کنند قرار دهیم آن گاه گوسفند، سگ، گربه، مرغ دریایی، افعی، گنجشک و مارمولک (حیواناتی که بیرون آب زندگی می کنند) کلاستر اول و ماهی قرمز، شاه ماهی و کوسه (حیواناتی که در آب زندگی می کنند) کلاستر دوم را تشکیل می دهند و قورباغه که می تواند هم در آب و هم در خشکی زندگی کند کلاستر سوم را تشکیل می دهد (C).



(c)

اگر وجودستون فقرات ملاک دسته بندی باشد تمام حیوانات در یک دسته قرار می گیرند.

ما می توانیم از ملاک دسته بندی مرکب استفاده کنیم. برای مثال اگر ملاک دسته بندی نحوه بدنی آوردن فرزندان و وجود شش ها باشد ما چهار نوع کلاستر داریم (d)



(d)

این مثال نشان می دهد که فرایند نسبت دادن اشیا به کلاسترها ممکن است به نتایج بسیار متفاوتی منجر شود. کلاسترینگ یکی از ابتدایی ترین فعالیت های ذهنی است که برای کنترل کردن مقادیر زیاد اطلاعات دریافت شده استفاده می شود. پردازش هر بخش از اطلاعات به عنوان یک موجودیت تک امکان پذیر نیست. بنابراین انسانها به دسته بندی موجودیت ها (حوادث، انسانها، اشیا و غیره) در کلاسترها روی می آورند

Let us now try to give some definitions for “clustering,” which, although they may not be universal, give us an idea of what clustering is. Let X be our data set, that is,

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}.$$

We define as an m -clustering of X , \mathfrak{R} , the partition of X into m sets (*clusters*), C_1, \dots, C_m , so that the following three conditions are met:

- $C_i \neq \emptyset, i = 1, \dots, m$
- $\bigcup_{i=1}^m C_i = X$
- $C_i \cap C_j = \emptyset, i \neq j, i, j = 1, \dots, m$

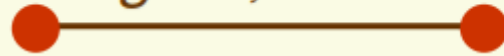
In addition, the vectors contained in a cluster C_i are “more similar” to each other and “less similar” to the feature vectors of the other clusters. Quantifying the terms *similar* and *dissimilar* depends very much on the types of clusters involved.

Proximity measure, either

similarity measure $s(\mathbf{x}_i, \mathbf{x}_k)$: large if $\mathbf{x}_i, \mathbf{x}_k$ are similar

dissimilarity(or distance) measure $d(\mathbf{x}_i, \mathbf{x}_k)$: small if $\mathbf{x}_i, \mathbf{x}_k$ are similar

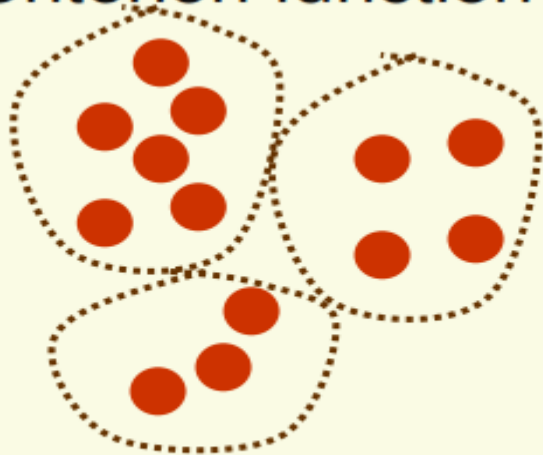
large d , small s



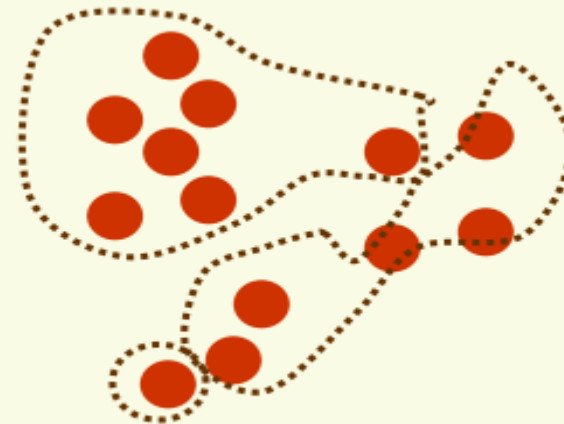
large s , small d



Criterion function to evaluate a clustering



good clustering



bad clustering

PROXIMITY MEASURES

A dissimilarity measure (DM) d on X is a function.

$$d : X \times X \rightarrow \mathcal{R}$$

where \mathcal{R} is the set of real numbers, such that

$$\exists d_0 \in \mathcal{R} : -\infty < d_0 \leq d(\mathbf{x}, \mathbf{y}) < +\infty, \quad \forall \mathbf{x}, \mathbf{y} \in X \quad (11.4)$$

$$d(\mathbf{x}, \mathbf{x}) = d_0, \quad \forall \mathbf{x} \in X \quad (11.5)$$

and

$$d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in X \quad (11.6)$$

If in addition

$$d(\mathbf{x}, \mathbf{y}) = d_0 \quad \text{if and only if} \quad \mathbf{x} = \mathbf{y} \quad (11.7)$$

and

$$d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}), \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X \quad (11.8)$$

d is called a *metric DM*. Inequality (11.8) is also known as the *triangular inequality*. Finally, equivalence (11.7) indicates that the minimum possible dissimilarity level value d_0 between any two vectors in X is achieved when they are identical. Sometimes we will refer to the dissimilarity level as distance, where the term is not used in its strict mathematical sense.

A *similarity measure* (SM) s on X is defined as

$$s : X \times X \rightarrow \mathcal{R}$$

such that

$$\exists s_0 \in \mathcal{R} : -\infty < s(\mathbf{x}, \mathbf{y}) \leq s_0 < +\infty, \quad \forall \mathbf{x}, \mathbf{y} \in X$$

$$s(\mathbf{x}, \mathbf{x}) = s_0, \quad \forall \mathbf{x} \in X$$

and

$$s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in X \quad (11.11)$$

If in addition

$$s(\mathbf{x}, \mathbf{y}) = s_0 \quad \text{if and only if} \quad \mathbf{x} = \mathbf{y} \quad (11.12)$$

and

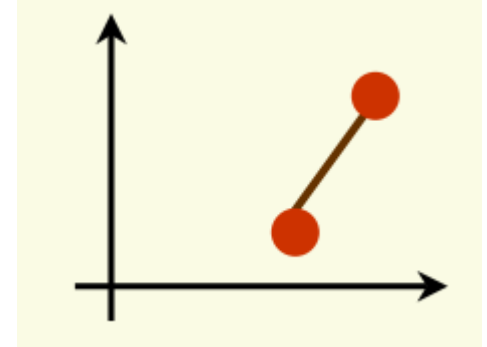
$$s(\mathbf{x}, \mathbf{y})s(\mathbf{y}, \mathbf{z}) \leq [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]s(\mathbf{x}, \mathbf{z}), \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X \quad (11.13)$$

s is called a *metric SM*.

Example 11.2

Let us consider the well-known Euclidean distance, d_2

$$d_2(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^l (x_i - y_i)^2}$$



where $\mathbf{x}, \mathbf{y} \in X$ and x_i, y_i are the i th coordinates of \mathbf{x} and \mathbf{y} , respectively. This is a dissimilarity measure on X , with $d_0 = 0$; that is, the minimum possible distance between two vectors of X is 0. Moreover, the distance of a vector from itself is equal to 0. Also, it is easy to observe that $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$.

In the sequel, we extend the

preceding definitions in order to measure “proximity” between subsets of X . Let U be a set containing subsets of X . That is, $D_i \subset X, i = 1, \dots, k$, and $U = \{D_1, \dots, D_k\}$. A *proximity measure* \wp on U is a function

$$\wp : U \times U \rightarrow \mathcal{R}$$

Equations (11.4)-(11.8) for dissimilarity measures and Eqs. (11.9)-(11.13) for similarity measures can now be repeated with D_i, D_j in the place of \mathbf{x} and \mathbf{y} and U in the place of X .

Usually, the proximity measures between two sets D_i and D_j are defined in terms of proximity measures between elements of D_i and D_j .

Example 11.3

Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6\}$ and $U = \{\{\mathbf{x}_1, \mathbf{x}_2\}, \{\mathbf{x}_1, \mathbf{x}_4\}, \{\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}, \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}\}$. Let us define the following dissimilarity function:

$$d_{\min}^{ss}(D_i, D_j) = \min_{\mathbf{x} \in D_i, \mathbf{y} \in D_j} d_2(\mathbf{x}, \mathbf{y})$$

where d_2 is the Euclidean distance between two vectors and $D_i, D_j \in U$.

The minimum possible value of d_{\min}^{ss} is $d_{\min,0}^{ss} = 0$. Also, $d_{\min}^{ss}(D_i, D_i) = 0$, since the Euclidean distance between a vector in D_i and itself is 0. In addition, it is easy to see that the commutative property holds. Thus, this dissimilarity function is a measure. It is not difficult to see that d_{\min}^{ss} is not a metric. Indeed, Eq. (11.7) for subsets of X does not hold in general, since the two sets D_i and D_j may have an element in common. Consider, for example the two sets $\{\mathbf{x}_1, \mathbf{x}_2\}$ and $\{\mathbf{x}_1, \mathbf{x}_4\}$ of U . Although they are different, their distance d_{\min}^{ss} is 0, since they both contain \mathbf{x}_1 .

Proximity Measures between Two Points

Real-Valued Vectors

A. Dissimilarity Measures

The most common DMs between real-valued vectors used in practice are:

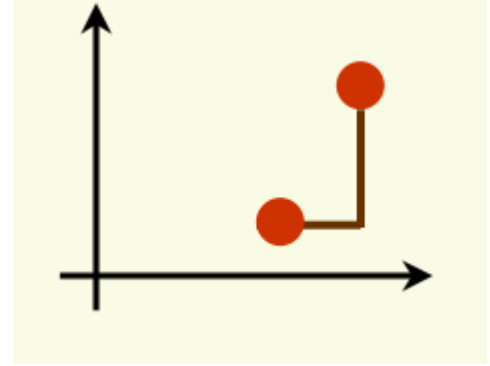
- The *weighted* l_p metric DMs, that is,

$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^l w_i |x_i - y_i|^p \right)^{1/p} \quad (11.14)$$

where x_i, y_i are the i th coordinates of \mathbf{x} and \mathbf{y} , $i = 1, \dots, l$, and $w_i \geq 0$ is the i th *weight coefficient*. They are used mainly on real-valued vectors. If $w_i = 1, i = 1, \dots, l$, we obtain the *unweighted* l_p metric DMs.

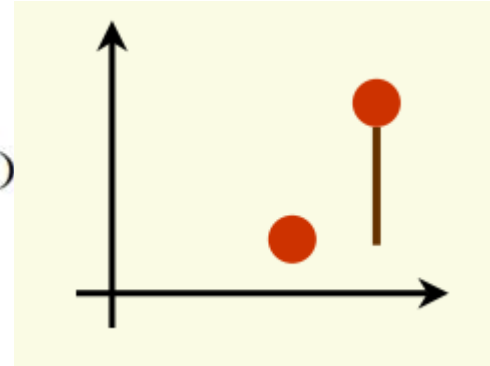
Special l_p metric DMs that are also encountered in practice are the (weighted) l_1 or *Manhattan norm*,

$$d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^l w_i |x_i - y_i| \quad (11.16)$$



and the (weighted) l_∞ norm,

$$d_\infty(\mathbf{x}, \mathbf{y}) = \max_{1 \leq i \leq l} w_i |x_i - y_i| \quad (11.17)$$



The l_1 and l_∞ norms may be viewed as overestimation and underestimation of the l_2 norm, respectively. Indeed, it can be shown that $d_\infty(\mathbf{x}, \mathbf{y}) \leq d_2(\mathbf{x}, \mathbf{y}) \leq d_1(\mathbf{x}, \mathbf{y})$ (see Problem 11.6). When $l = 1$ all l_p norms coincide.

Based on these DMs, we can define corresponding SMs as

$$s_p(\mathbf{x}, \mathbf{y}) = b_{\max} - d_p(\mathbf{x}, \mathbf{y}).$$

- Some additional DMs are the following

$$d_G(\mathbf{x}, \mathbf{y}) = -\log_{10} \left(1 - \frac{1}{l} \sum_{j=1}^l \frac{|x_j - y_j|}{b_j - a_j} \right) \quad (11.18)$$

where b_j and a_j are the maximum and the minimum values among the j th features of the N vectors of X , respectively. It can easily be shown that $d_G(\mathbf{x}, \mathbf{y})$ is a metric DM. Notice that the value of $d_G(\mathbf{x}, \mathbf{y})$ depends not only on \mathbf{x} and \mathbf{y} but also on the whole of X . Thus, if $d_G(\mathbf{x}, \mathbf{y})$ is the distance between two vectors \mathbf{x} and \mathbf{y} that belong to a set X and $d'_G(\mathbf{x}, \mathbf{y})$ is the distance between the same two vectors when they belong to a different set X' , then, in general,

$$d_G(\mathbf{x}, \mathbf{y}) \neq d'_G(\mathbf{x}, \mathbf{y}).$$

Example 11.4

Consider the three-dimensional vectors $\mathbf{x} = [0, 1, 2]^T$, $\mathbf{y} = [4, 3, 2]^T$. Then, assuming that all w_i 's are equal to 1, $d_1(\mathbf{x}, \mathbf{y}) = 6$, $d_2(\mathbf{x}, \mathbf{y}) = 2\sqrt{5}$, and $d_\infty(\mathbf{x}, \mathbf{y}) = 4$. Notice that $d_\infty(\mathbf{x}, \mathbf{y}) < d_2(\mathbf{x}, \mathbf{y}) < d_1(\mathbf{x}, \mathbf{y})$.

Assume now that these vectors belong to a data set X that contains N vectors with maximum values per feature 10, 12, 13 and minimum values per feature 0, 0.5, 1, respectively. Then $d_G(\mathbf{x}, \mathbf{y}) = 0.0922$. If, on the other hand, \mathbf{x} and \mathbf{y} belong to an X' with the maximum (minimum) values per feature being 20, 22, 23 ($-10, -9.5, -9$), respectively, then $d_G(\mathbf{x}, \mathbf{y}) = 0.0295$.

Finally, $d_Q(\mathbf{x}, \mathbf{y}) = 0.6455$.

$$d_Q(\mathbf{x}, \mathbf{y}) = \sqrt{\frac{1}{l} \sum_{j=1}^l \left(\frac{x_j - y_j}{x_j + y_j} \right)^2}$$

$$d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^l w_i |x_i - y_i|$$

$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^l w_i |x_i - y_i|^p \right)^{1/p}$$

$$d_\infty(\mathbf{x}, \mathbf{y}) = \max_{1 \leq i \leq l} w_i |x_i - y_i|$$

$$d_G(\mathbf{x}, \mathbf{y}) = -\log_{10} \left(1 - \frac{1}{l} \sum_{j=1}^l \frac{|x_j - y_j|}{b_j - a_j} \right)$$

B. Similarity Measures

The most common similarity measures for real-valued vectors used in practice are:

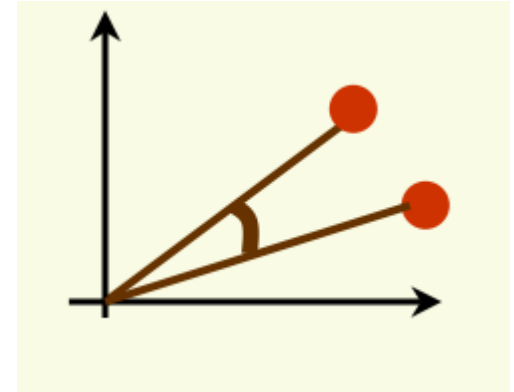
The inner product. It is defined as

$$s_{\text{inner}}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^l x_i y_i.$$

In most cases, the inner product is used when the vectors \mathbf{x} and \mathbf{y} are normalized, so that they have the same length a . In these cases, the upper and the lower bounds of s_{inner} are $+a^2$ and $-a^2$, respectively, and $s_{\text{inner}}(\mathbf{x}, \mathbf{y})$ depends exclusively on the angle between \mathbf{x} and \mathbf{y} .

- cosine similarity measure:

$$s_{\text{cosine}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$



where $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^l x_i^2}$ and $\|\mathbf{y}\| = \sqrt{\sum_{i=1}^l y_i^2}$ are the lengths of the vectors \mathbf{x} and \mathbf{y} , respectively. This measure is invariant to rotations but not to linear transformations.

- Pearson's correlation coefficient.

This measure can be expressed as

$$r_{\text{Pearson}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}_d^T \mathbf{y}_d}{\|\mathbf{x}_d\| \|\mathbf{y}_d\|}$$

where $\mathbf{x}_d = [x_1 - \bar{x}, \dots, x_l - \bar{x}]^T$ and $\mathbf{y}_d = [y_1 - \bar{y}, \dots, y_l - \bar{y}]^T$, with x_i, y_i being the i th coordinates of \mathbf{x} and \mathbf{y} , respectively, and $\bar{x} = \frac{1}{l} \sum_{i=1}^l x_i$, $\bar{y} = \frac{1}{l} \sum_{i=1}^l y_i$. Usually, \mathbf{x}_d and \mathbf{y}_d are called difference vectors. Clearly, $r_{\text{Pearson}}(\mathbf{x}, \mathbf{y})$ takes values between -1 and $+1$.

A related dissimilarity measure can be defined as

$$D(\mathbf{x}, \mathbf{y}) = \frac{1 - r_{\text{Pearson}}(\mathbf{x}, \mathbf{y})}{2}$$

- Another commonly used SM is the *Tanimoto measure*, which is also known as Tanimoto distance [Tani 58]. It may be used for real- as well as for discrete-valued vectors. It is defined as

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x}^T \mathbf{y}} \quad (11.23)$$

Proximity Functions between a Point and a Set

In many clustering schemes, a vector \mathbf{x} is assigned to a cluster C taking into account the proximity between \mathbf{x} and C , $\wp(\mathbf{x}, C)$.

■ The max proximity function:

$$\wp_{\max}^{ps}(\mathbf{x}, C) = \max_{\mathbf{y} \in C} \wp(\mathbf{x}, \mathbf{y})$$

■ The min proximity function:

$$\wp_{\min}^{ps}(\mathbf{x}, C) = \min_{\mathbf{y} \in C} \wp(\mathbf{x}, \mathbf{y})$$

■ The average proximity function:

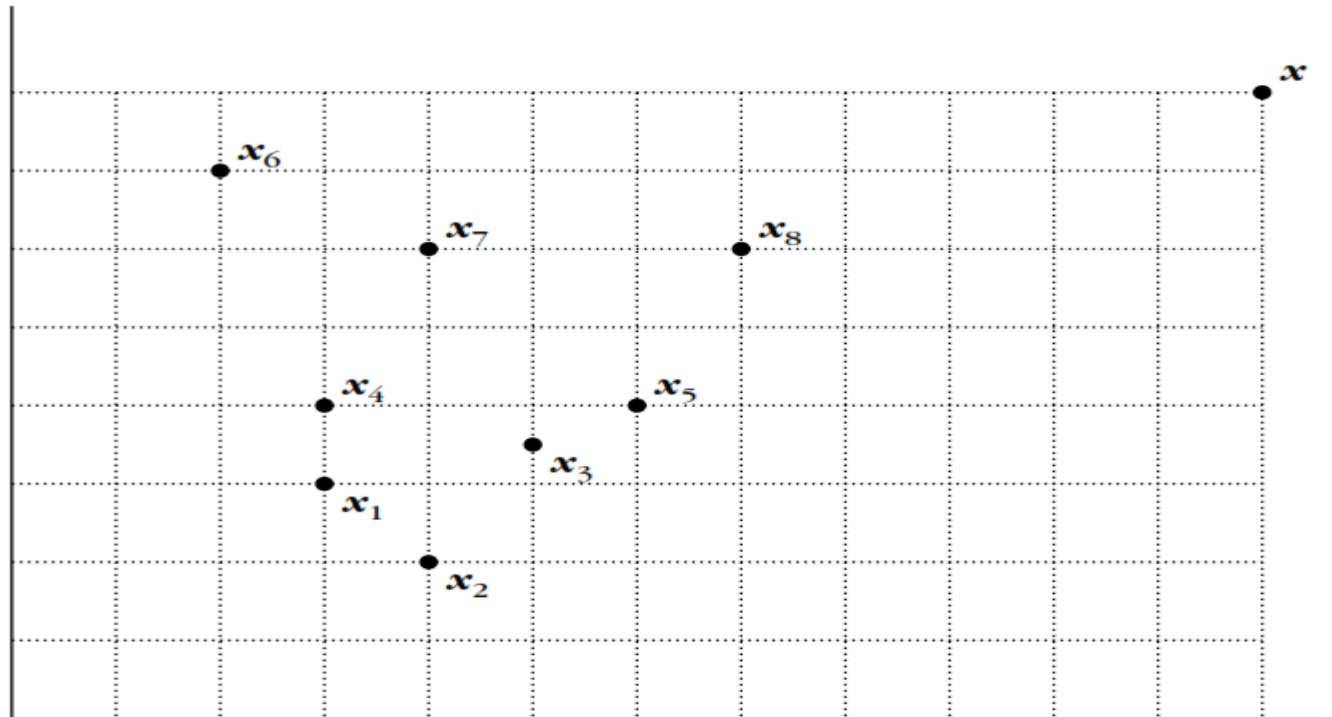
$$\wp_{\text{avg}}^{ps}(\mathbf{x}, C) = \frac{1}{n_C} \sum_{\mathbf{y} \in C} \wp(\mathbf{x}, \mathbf{y})$$

where n_C is the cardinality of C .

In these definitions, $\wp(\mathbf{x}, \mathbf{y})$ may be any proximity measure between two points.

Example 11.9

Let $C = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6, \mathbf{x}_7, \mathbf{x}_8\}$, where $\mathbf{x}_1 = [1.5, 1.5]^T$, $\mathbf{x}_2 = [2, 1]^T$, $\mathbf{x}_3 = [2.5, 1.75]^T$, $\mathbf{x}_4 = [1.5, 2]^T$, $\mathbf{x}_5 = [3, 2]^T$, $\mathbf{x}_6 = [1, 3.5]^T$, $\mathbf{x}_7 = [2, 3]^T$, $\mathbf{x}_8 = [3.5, 3]^T$, and let $\mathbf{x} = [6, 4]^T$ (see Figure 11.6). Assume that the Euclidean distance is used to measure the dissimilarity between two points. Then $d_{\max}^{ps}(\mathbf{x}, C) = \max_{\mathbf{y} \in C} d(\mathbf{x}, \mathbf{y}) = d(\mathbf{x}, \mathbf{x}_1) = 5.15$. For the other two distances we have $d_{\min}^{ps}(\mathbf{x}, C) = \min_{\mathbf{y} \in C} d(\mathbf{x}, \mathbf{y}) = d(\mathbf{x}, \mathbf{x}_8) = 2.69$ and $d_{\text{avg}}^{ps}(\mathbf{x}, C) = \frac{1}{n_C} \sum_{\mathbf{y} \in C} d(\mathbf{x}, \mathbf{y}) = \frac{1}{8} \sum_{i=1}^8 d(\mathbf{x}, \mathbf{x}_i) = 4.33$.



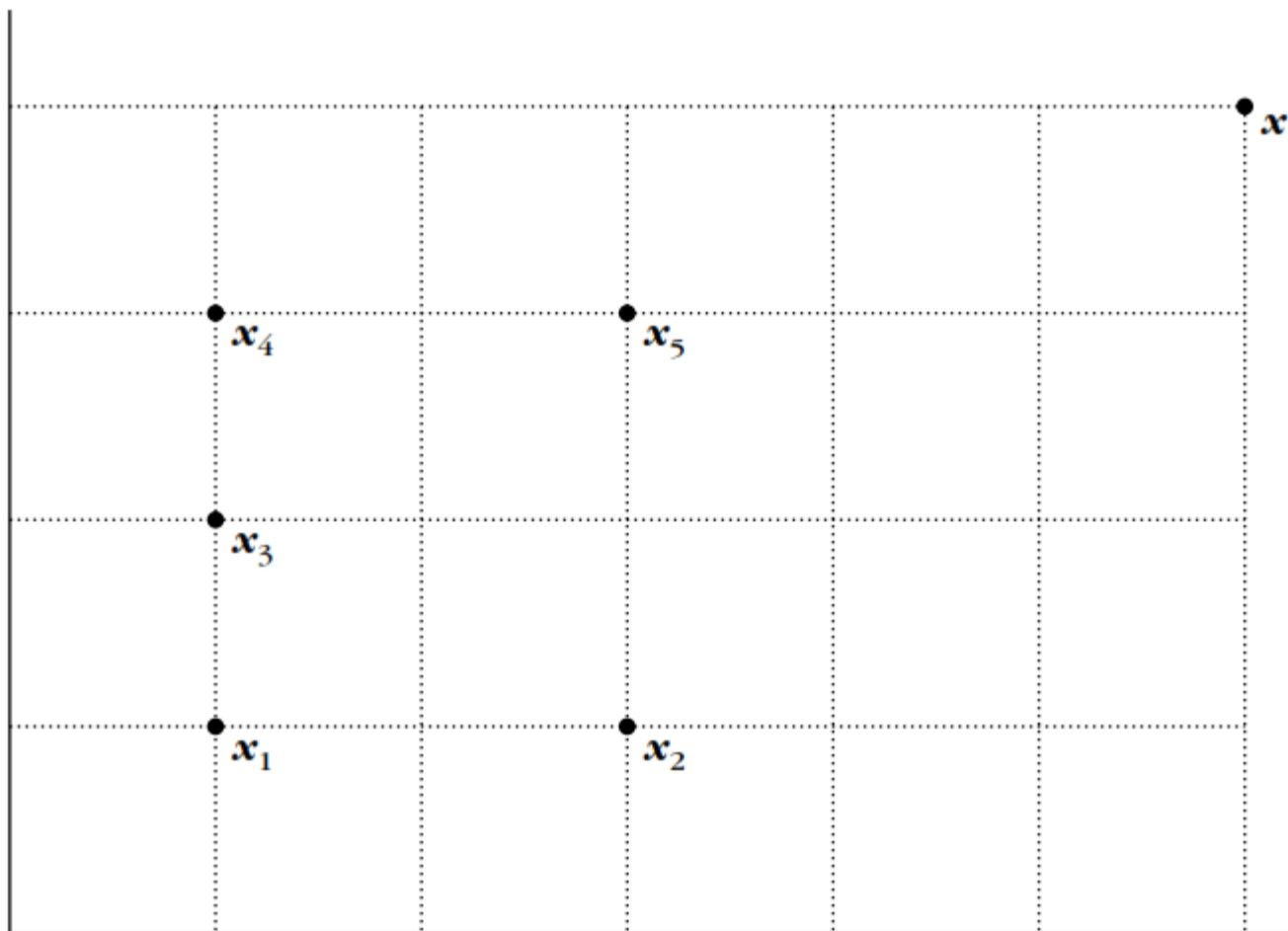
Typical choices for a point representative of a cluster are:

- The mean vector (or mean point)

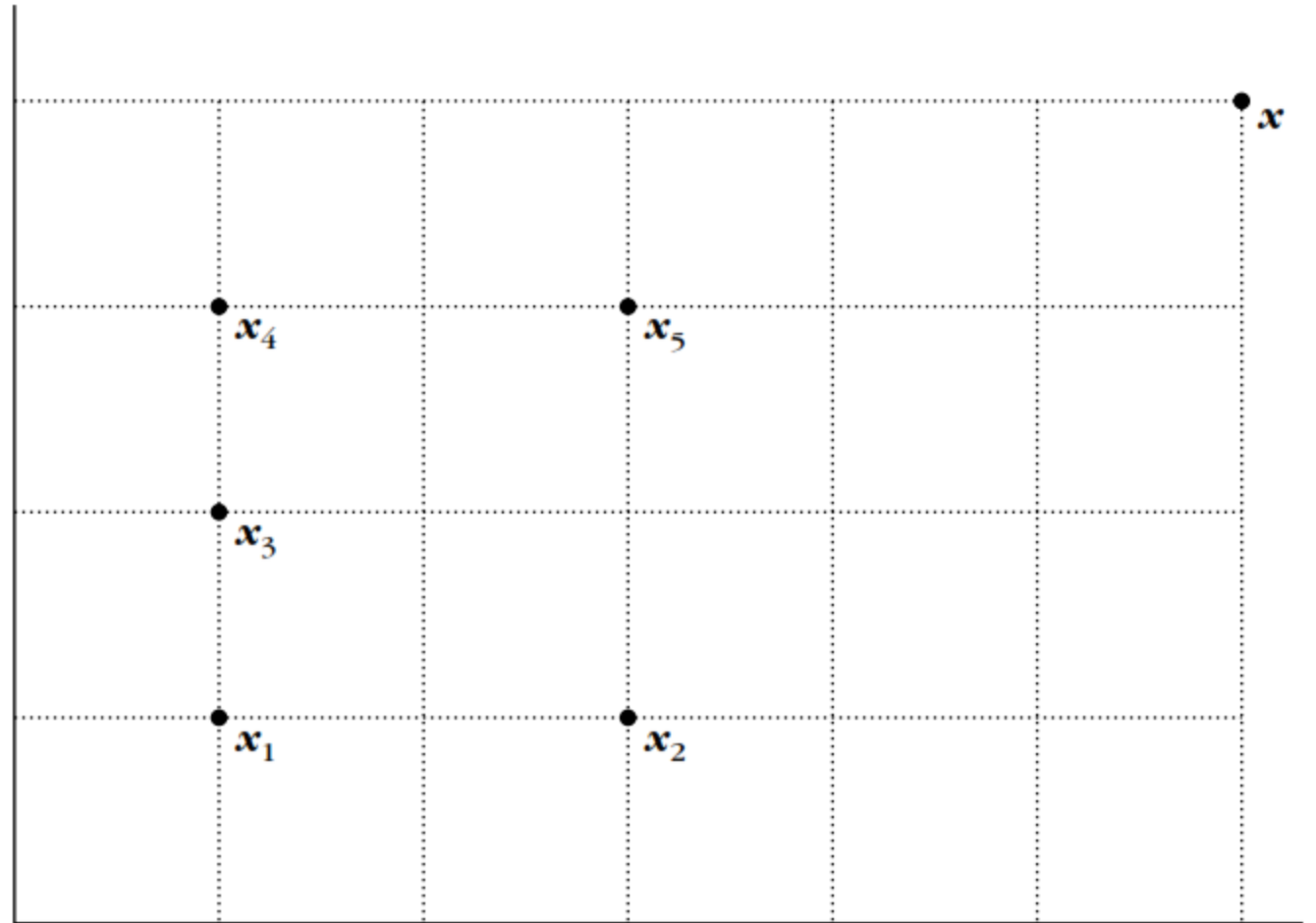
$$\mathbf{m}_p = \frac{1}{n_C} \sum_{\mathbf{y} \in C} \mathbf{y}$$

Example 11.10

Let $C = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$, where $\mathbf{x}_1 = [1, 1]^T$, $\mathbf{x}_2 = [3, 1]^T$, $\mathbf{x}_3 = [1, 2]^T$, $\mathbf{x}_4 = [1, 3]^T$, and $\mathbf{x}_5 = [3, 3]^T$ (see Figure 11.8). All points lie in the discrete space $\{0, 1, 2, \dots, 6\}^2$. We use the Euclidean distance to measure the dissimilarity between two vectors in C . The mean point of C is $\mathbf{m}_p = [1.8, 2]^T$. It is clear that \mathbf{m}_p lies outside the space where the elements of C belong.



To find the mean center \mathbf{m}_C , we compute, for each point $\mathbf{x}_i \in C$, $i = 1, \dots, 5$, the sum A_i of its distances from all other points of C . The resulting values are $A_1 = 7.83$, $A_2 = 9.06$, $A_3 = 6.47$, $A_4 = 7.83$, $A_5 = 9.06$. The minimum of these values is A_3 . Thus, \mathbf{x}_3 is the mean center of C .



11.3 PROBLEMS

11.1 Let s be a metric similarity measure on X with $s(\mathbf{x}, \mathbf{y}) > 0, \forall \mathbf{x}, \mathbf{y} \in X$ and $d(\mathbf{x}, \mathbf{y}) = a/s(\mathbf{x}, \mathbf{y})$, with $a > 0$. Prove that d is a metric dissimilarity measure.

11.2 Prove that the Euclidean distance satisfies the triangular inequality.

Hint: Use the Minkowski inequality, which states that for a positive integer p and two vectors $\mathbf{x} = [x_1, \dots, x_l]^T$ and $\mathbf{y} = [y_1, \dots, y_l]^T$ it holds that

$$\left(\sum_{i=1}^l |x_i + y_i|^p \right)^{1/p} \leq \left(\sum_{i=1}^l |x_i|^p \right)^{1/p} + \left(\sum_{i=1}^l |y_i|^p \right)^{1/p}$$

11.6 Prove that

$$d_\infty(\mathbf{x}, \mathbf{y}) \leq d_2(\mathbf{x}, \mathbf{y}) \leq d_1(\mathbf{x}, \mathbf{y})$$

for any two vectors \mathbf{x} and \mathbf{y} in X .

انواع کلاسترها

کلاسترها انواع مختلفی دارند که در ادامه به تعدادی از آنها اشاره خواهد شد:

✓ **کلاسترهای بخوبی جدا شده:**

مجموعه نقاط داخل این کلاستر نسبت به نقاط خارج آن به یکدیگر بسیار شبیه ترند.

✓ **کلاسترهای مبتنی به مرکز:**

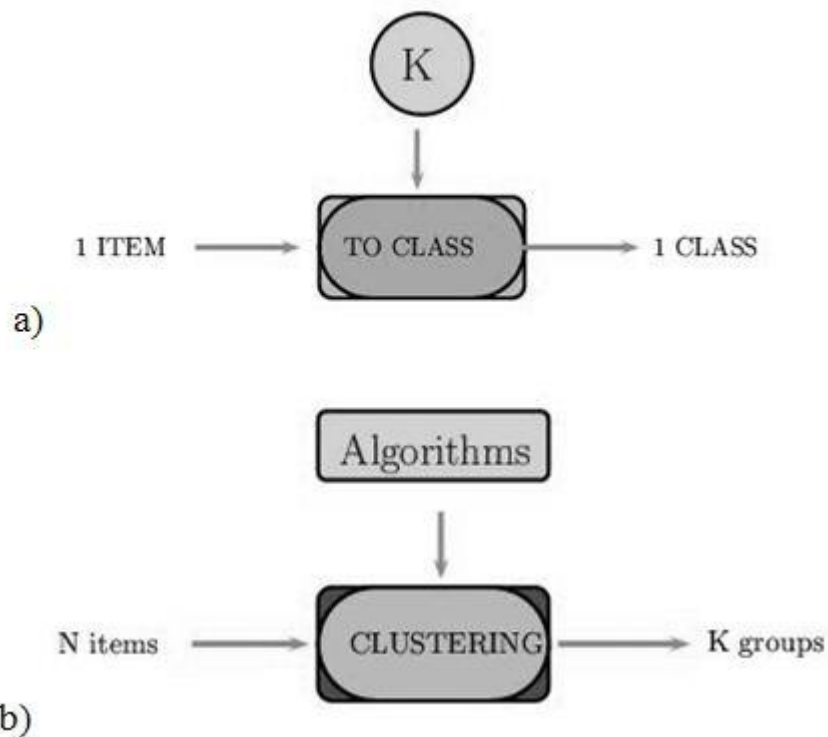
مجموعه نقاط داخل این کلاستر به مرکز کلاستر نسبت به مراکز کلاسترهای دیگر بسیار نزدیکترند.

▪ **کلاسترهای مبتنی بر مجاورت و نزدیکی:**

مجموعه نقاط داخل این کلاستر به یک یا تعداد بیشتری از نقاط داخل کلاستر نسبت به نقاط خارج آن شبیه اند.

خوشه‌بندی در مقابل طبقه‌بندی

در طبقه‌بندی، دسته‌ها (کلاس) از پیش مشخص شده‌اند و هر داده به یک دسته از پیش مشخص شده تخصیص می‌یابد. ولی در خوشه‌بندی هیچ اطلاعی از کلاسهای موجود درون داده‌ها وجود ندارد و به عبارتی خود خوشه‌ها نیز از داده‌ها استخراج می‌شوند. در شکل زیر تفاوت بین خوشه‌بندی و طبقه‌بندی بهتر نشان داده شده است.



(a) در طبقه‌بندی با استفاده یک سری اطلاعات اولیه داده‌ها به دسته‌های معلومی نسبت داده می‌شوند

(b) در خوشه‌بندی داده‌ها با توجه به الگوریتم انتخاب شده به خوشه‌هایی نسبت داده می‌شوند

کاربردها

یک جستجوی ساده در وب یا حتی در پایگاه داده ی یک کتابخانه، کاربرد شگفت انگیز خوشه بندی را برای ما آشکار می سازد. از آنجا که خوشه بندی یک روش یادگیری بدون نظارت محسوب می گردد، در موارد بسیاری می تواند کاربرد داشته باشد:

- **بازاریابی:** دسته بندی مشتری ها به دسته هایی بر حسب رفتارها و نیازهای آن ها از طریق مجموعه زیادی از ویژگی ها و آخرین خریدهای آن ها.
- **زیست شناسی:** دسته بندی حیوانات و گیاهان از روی ویژگی های آنها
- **کتابداری:** دسته بندی کتابها
- **نقشه برداری شهری:** دسته بندی خانه ها بر اساس نوع و موقعیت جغرافیایی آنها.
- **مطالعات زلزله نگاری:** تشخیص مناطق حادثه خیز بر اساس مشاهدات قبلی.
- **وب:** دسته بندی اسناد، دسته بندی مشتریان مراجعه کننده به سایتها و
- **داده کاوی:** کشف اطلاعات و ساختار جدید از داده های موجود
- **در تشخیص گفتار:** در ساخت کتاب کد از بردارهای ویژگی، در تقسیم کردن گفتار بر حسب گویندگان آن و یا فشرده سازی گفتار
- **در تقسیم بندی تصاویر:** تقسیم بندی تصاویر پزشکی و یا ماهواره ای
- **بیمه:** تشخیص افراد متقلب، تشخیص افرادی که بیمه موتور دارند و بیشترین میزان درخواست از بیمه را نیز در سال مشخصی داشته اند.

ایرادات

متأسفانه چندین مسئله در خصوص روش‌های خوشه‌بندی مطرح است که هنوز به شکل کامل پاسخ داده نشده‌اند. و همچنان تلاش‌های بسیاری به منظور حل آنها انجام گرفته است:

❖ روش‌های خوشه‌بندی قادر نیستند تمامی نیازهای مسائل را به طور هم‌زمان برآورده کنند.

❖ به دلیل پیچیدگی محاسباتی زیاد در برخورد با مجموعه داده‌های بزرگ (تعداد داده زیاد و تعداد ویژگی‌های زیاد برای هر داده) عملی نیستند.

❖ به دلیل وابستگی شدید به تعریف معیار شباهت بین داده‌ها، در مسائلی که تعریف معیار شباهت مشکل باشد نتایج مطلوبی تولید نمی‌کنند. (در داده‌ها با تعداد ویژگی زیاد)

AGGLOMERATIVE ALGORITHMS

Let $g(C_i, C_j)$ be a function defined for all possible pairs of clusters of X . This function measures the proximity between C_i and C_j . Let t denote the current level of hierarchy. Then, the general agglomerative scheme may be stated as follows:

Generalized Agglomerative Scheme (GAS)

■ Initialization:

- Choose $\mathfrak{R}_0 = \{C_i = \{\mathbf{x}_i\}, i = 1, \dots, N\}$ as the initial clustering.
- $t = 0$.

■ Repeat:

- $t = t + 1$
- Among all possible pairs of clusters (C_r, C_s) in \mathfrak{R}_{t-1} find the one, say (C_i, C_j) , such that

$$g(C_i, C_j) = \begin{cases} \min_{r,s} g(C_r, C_s), & \text{if } g \text{ is a dissimilarity function} \\ \max_{r,s} g(C_r, C_s), & \text{if } g \text{ is a similarity function} \end{cases} \quad (13.1)$$

- Define $C_q = C_i \cup C_j$ and produce the new clustering $\mathfrak{R}_t = (\mathfrak{R}_{t-1} - \{C_i, C_j\}) \cup \{C_q\}$.

■ Until all vectors lie in a single cluster.

Example 13.1

Let $X = \{\mathbf{x}_i, i = 1, \dots, 5\}$, with $\mathbf{x}_1 = [1, 1]^T$, $\mathbf{x}_2 = [2, 1]^T$, $\mathbf{x}_3 = [5, 4]^T$, $\mathbf{x}_4 = [6, 5]^T$, and $\mathbf{x}_5 = [6.5, 6]^T$. The pattern matrix of X is

$$D(X) = \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 5 & 4 \\ 6 & 5 \\ 6.5 & 6 \end{bmatrix}$$

and its corresponding dissimilarity matrix, when the Euclidean distance is in use, is

$$P(X) = \begin{bmatrix} 0 & 1 & 5 & 6.4 & 7.4 \\ 1 & 0 & 4.2 & 5.7 & 6.7 \\ 5 & 4.2 & 0 & 1.4 & 2.5 \\ 6.4 & 5.7 & 1.4 & 0 & 1.1 \\ 7.4 & 6.7 & 2.5 & 1.1 & 0 \end{bmatrix}$$

When the Tanimoto measure is used, the similarity matrix of X becomes

$$P'(X) = \begin{bmatrix} 1 & 0.75 & 0.26 & 0.21 & 0.18 \\ 0.75 & 1 & 0.44 & 0.35 & 0.20 \\ 0.26 & 0.44 & 1 & 0.96 & 0.90 \\ 0.21 & 0.35 & 0.96 & 1 & 0.98 \\ 0.18 & 0.20 & 0.90 & 0.98 & 1 \end{bmatrix}$$

Note that in $P(X)$ all diagonal elements are 0, since $d_2(\mathbf{x}, \mathbf{x}) = 0$, while in $P'(X)$ all diagonal elements are equal to 1, since $s_T(\mathbf{x}, \mathbf{x}) = 1$.

$$g(C_i, C_j) = d_{\min}^{SS}(C_i, C_j)$$

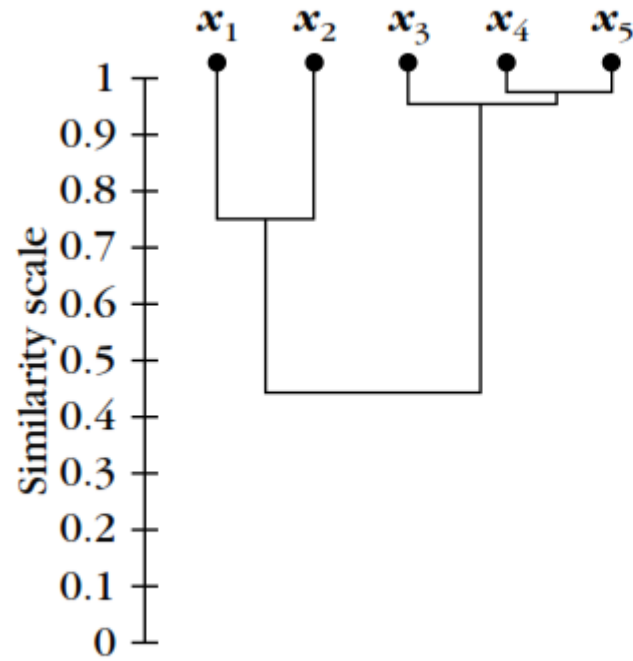
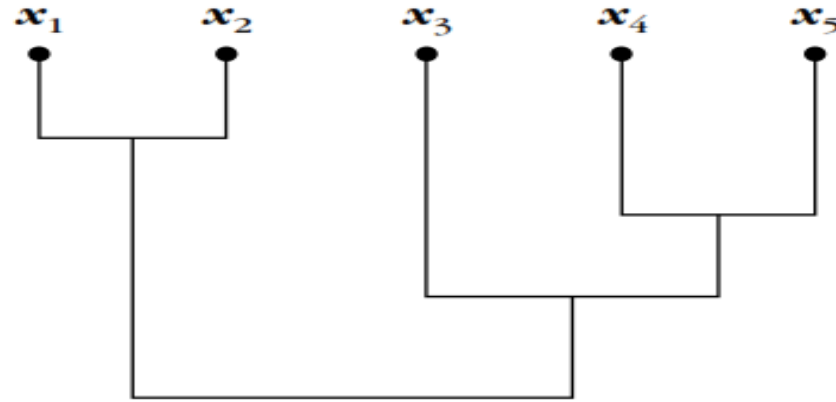
$\{\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}\}$

$\{\{x_1, x_2\}, \{x_3\}, \{x_4\}, \{x_5\}\}$

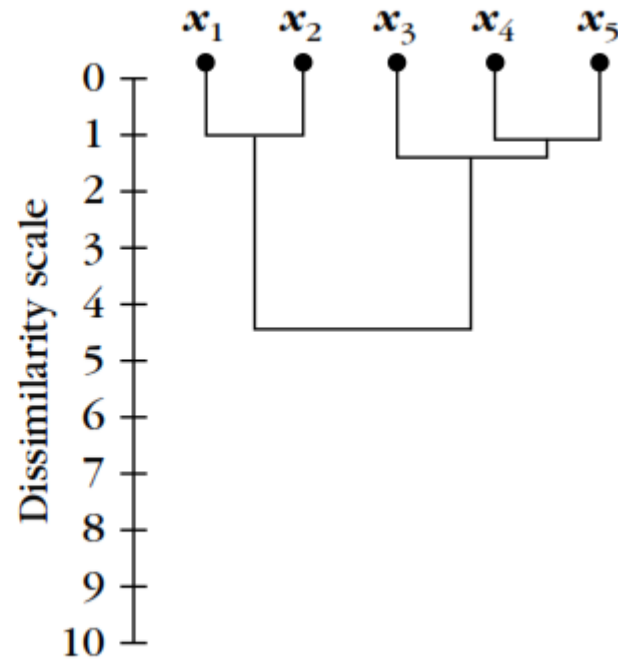
$\{\{x_1, x_2\}, \{x_3\}, \{x_4, x_5\}\}$

$\{\{x_1, x_2\}, \{x_3, x_4, x_5\}\}$

$\{\{x_1, x_2, x_3, x_4, x_5\}\}$

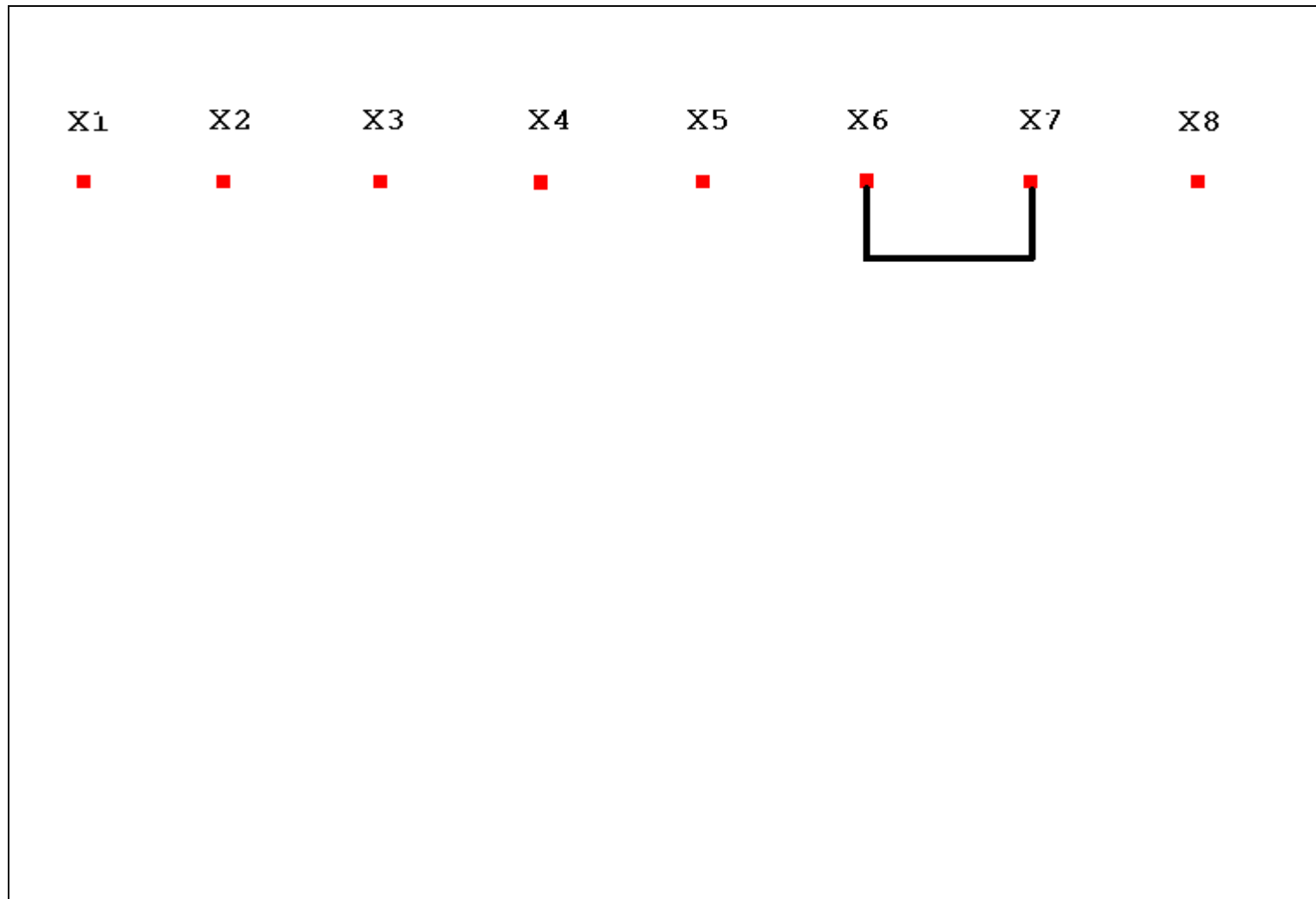


(a)

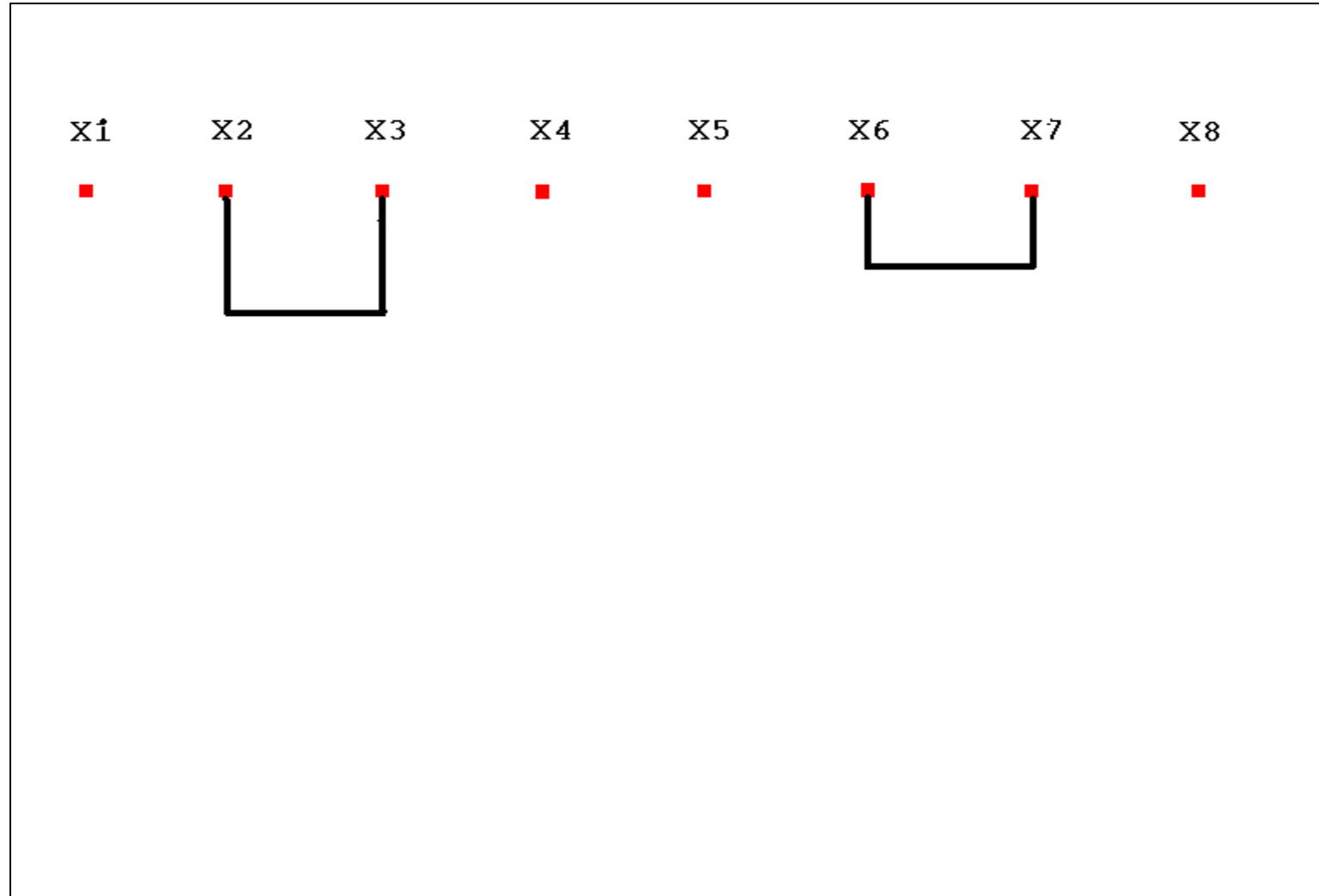


(b)

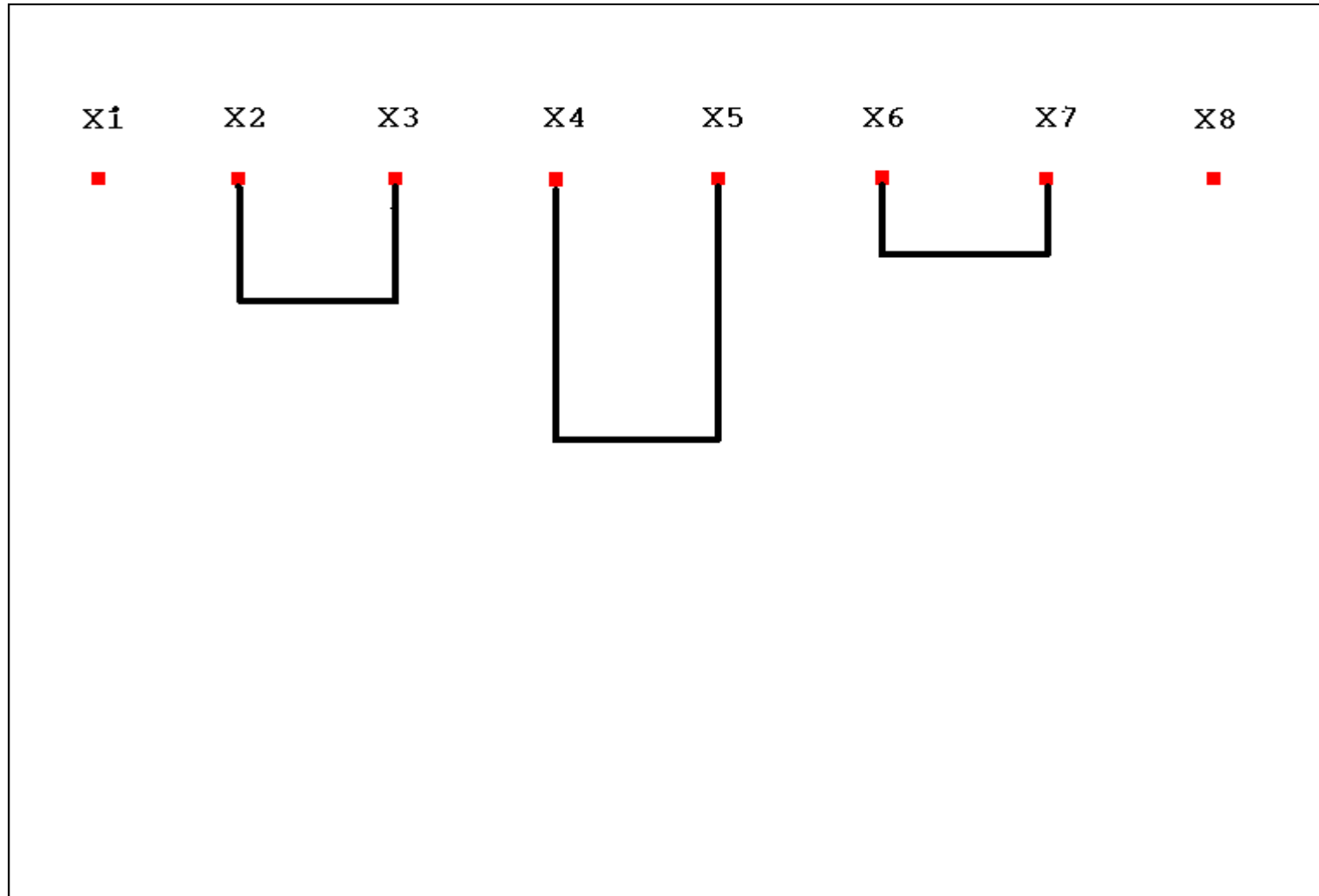
Nearest Neighbor, Level 2, k = 7 clusters.



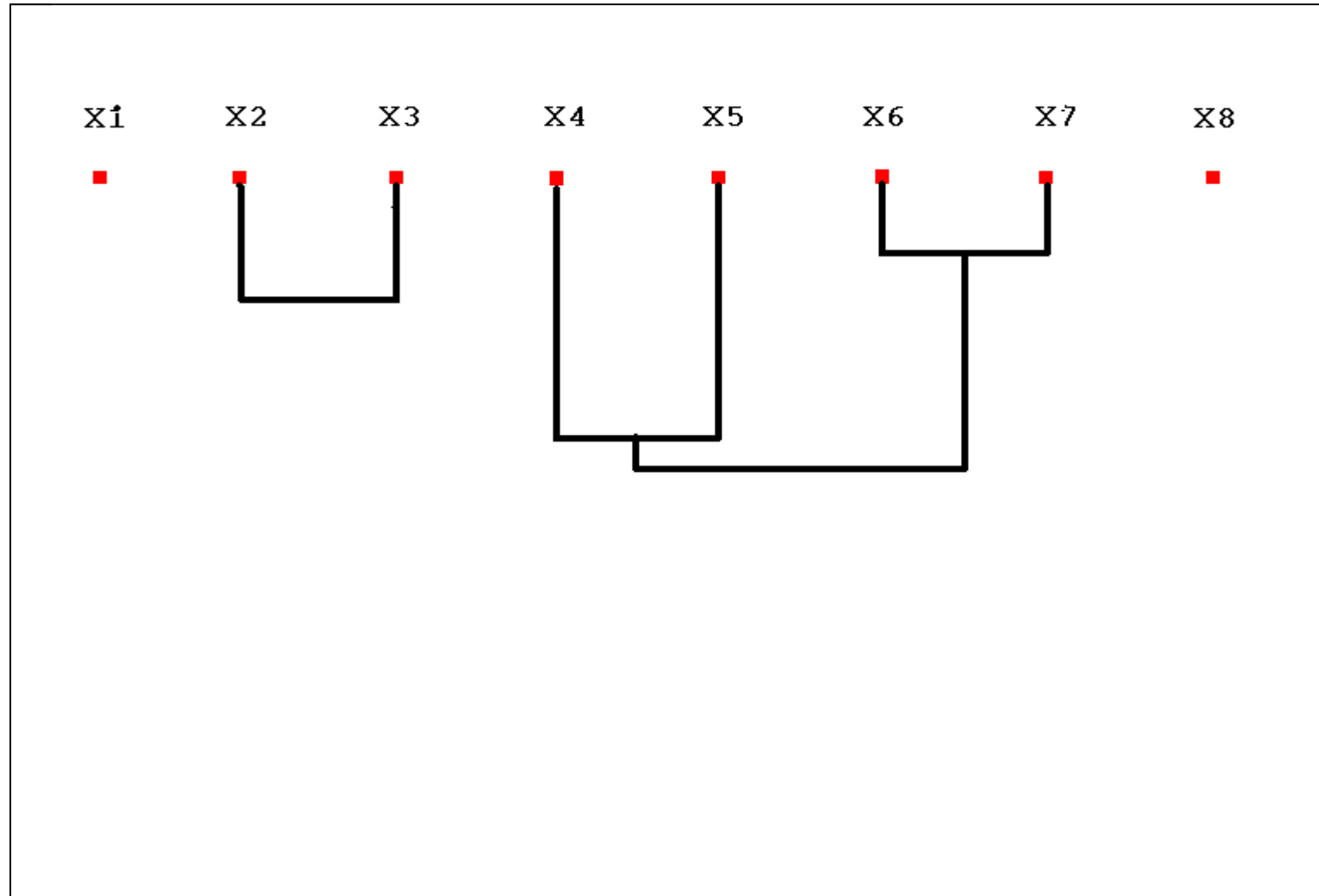
Nearest Neighbor, Level 3, k = 6 clusters.



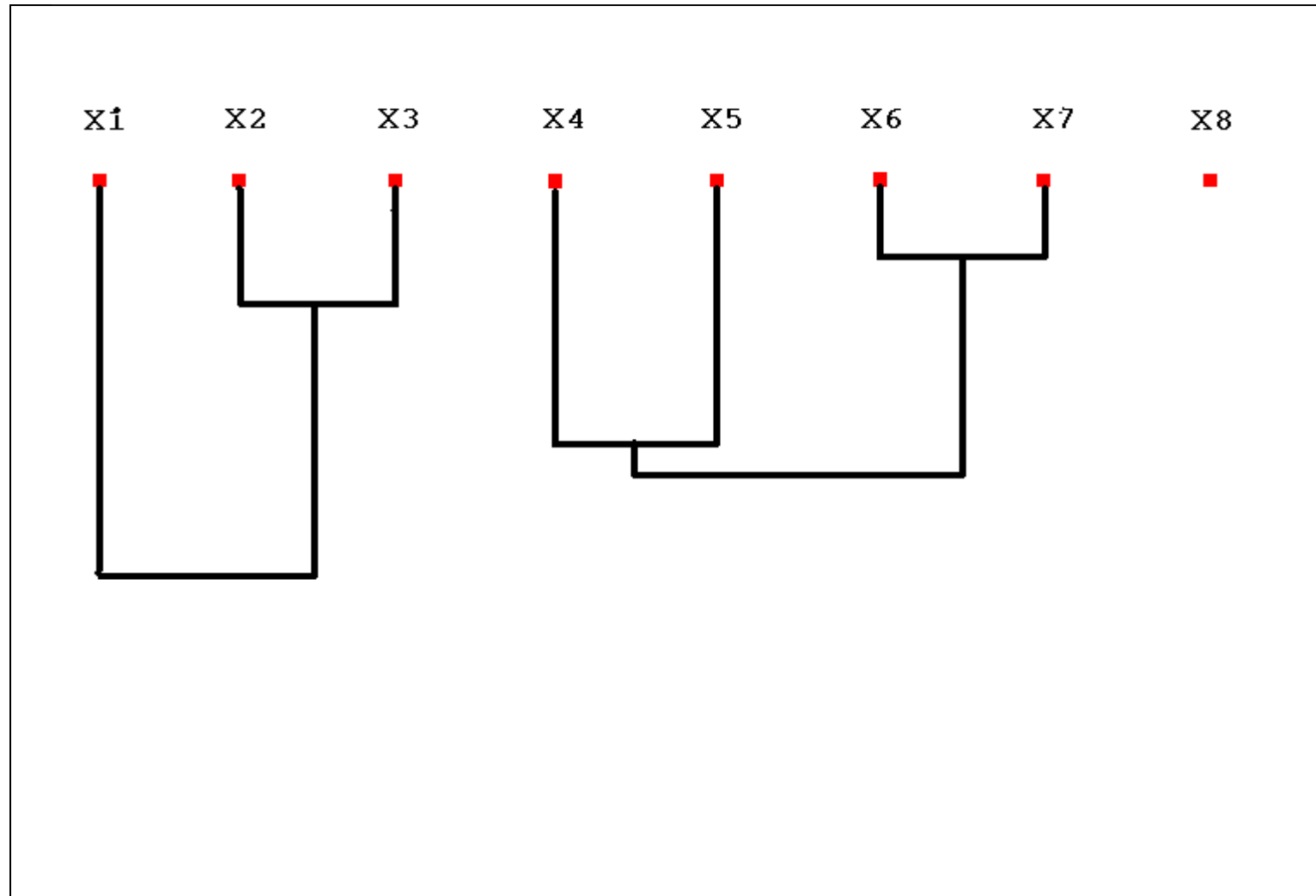
Nearest Neighbor, Level 4, k = 5 clusters.



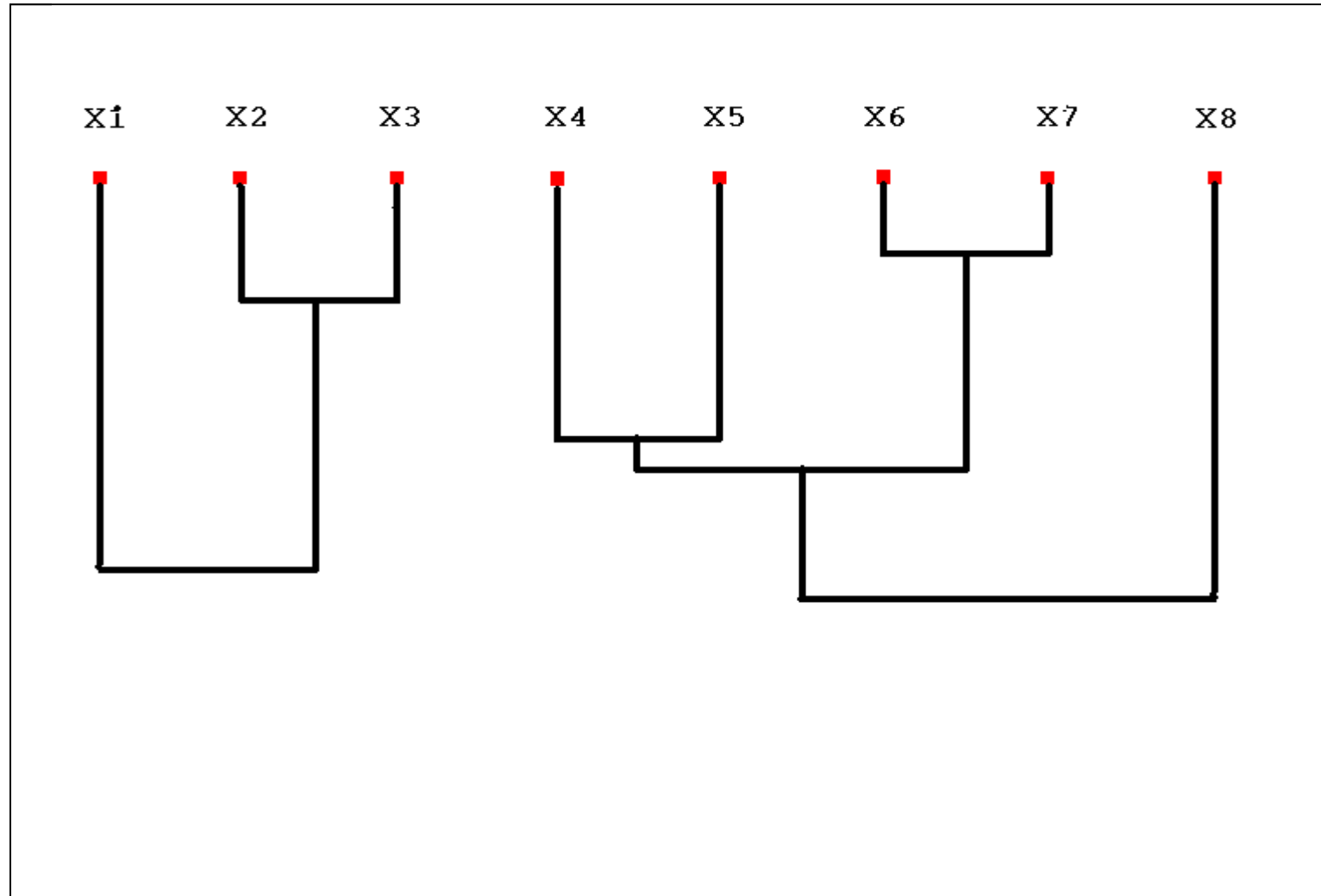
Nearest Neighbor, Level 5, k = 4 clusters.



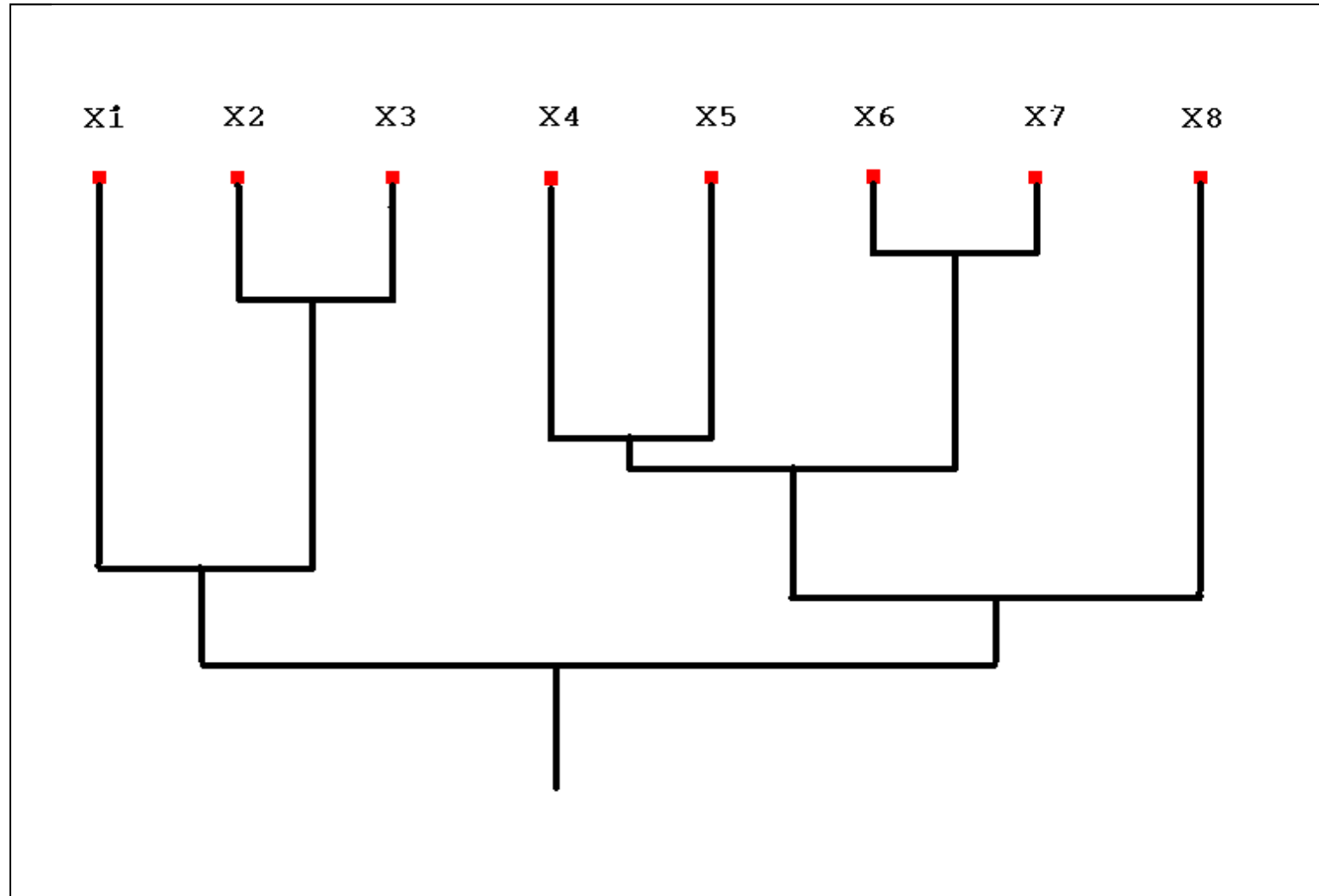
Nearest Neighbor, Level 6, $k = 3$ clusters.

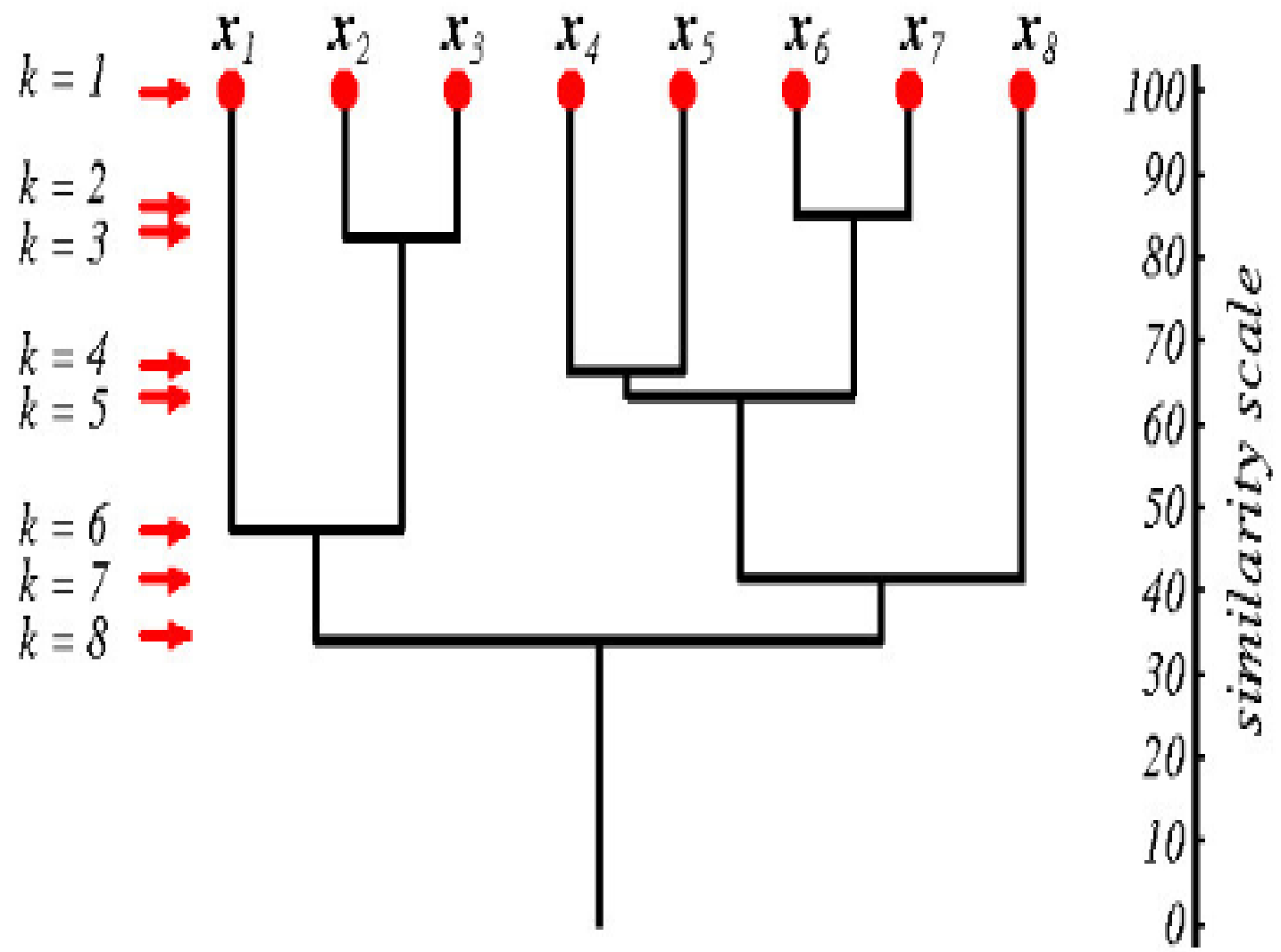


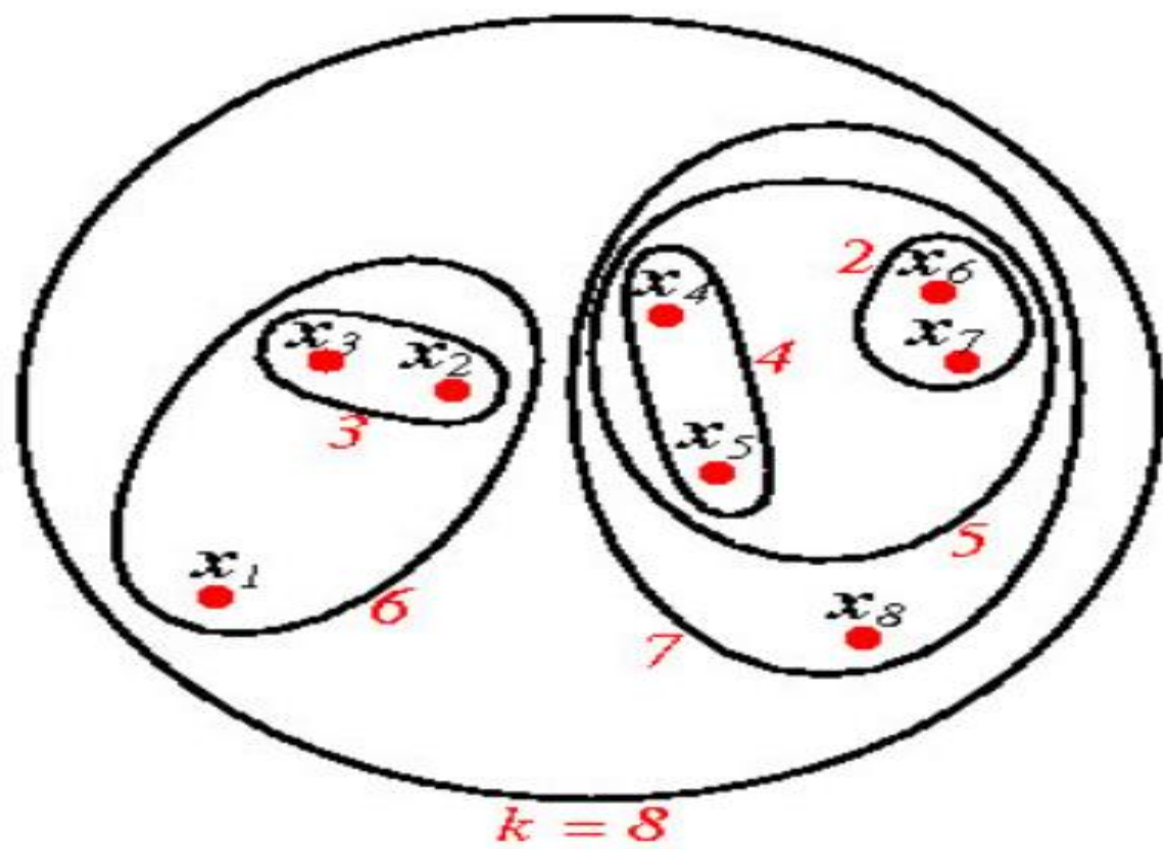
Nearest Neighbor, Level 7, k = 2 clusters.



Nearest Neighbor, Level 8, k = 1 cluster.







K- Means

این روش علی‌رغم سادگی آن یک روش پایه برای بسیاری از روش‌های خوشه‌بندی دیگر (مانند خوشه‌بندی فازی) محسوب می‌شود. این روش روشی انحصاری و مسطح محسوب می‌شود. برای این الگوریتم شکل‌های مختلفی بیان شده است. ولی همه آنها دارای روالی تکراری هستند که برای تعدادی ثابت از خوشه‌ها سعی در تخمین موارد زیر دارند:

- بدست آوردن نقاطی به عنوان مراکز خوشه‌ها. این نقاط در واقع همان میانگین نقاط متعلق به هر خوشه هستند.
 - نسبت دادن هر نمونه داده به یک خوشه که آن داده کمترین فاصله تا مرکز آن خوشه را دارا باشد.
- در نوع ساده‌ای از این روش ابتدا به تعداد خوشه‌های مورد نیاز نقاطی به صورت تصادفی انتخاب می‌شود. سپس در داده‌ها با توجه با میزان نزدیکی (شبهت) به یکی از این خوشه‌ها نسبت داده می‌شوند و بدین ترتیب خوشه‌های جدیدی حاصل می‌شود. با تکرار همین روال می‌توان در هر تکرار با میانگین‌گیری از داده‌ها مراکز جدیدی برای آنها محاسبه کرد و مجدداً داده‌ها را به خوشه‌های جدید نسبت داد. این روند تا زمانی ادامه پیدا می‌کند که دیگر تغییری در داده‌ها حاصل نشود. تابع زیر به عنوان تابع هدف مطرح است.

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

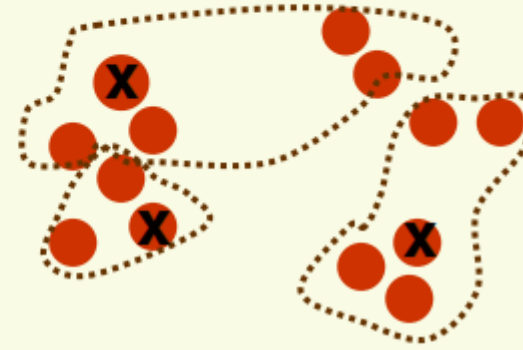
که $\| \|$ معیار فاصله بین نقاط و c_j مرکز خوشه j ام است.

K-means Clustering

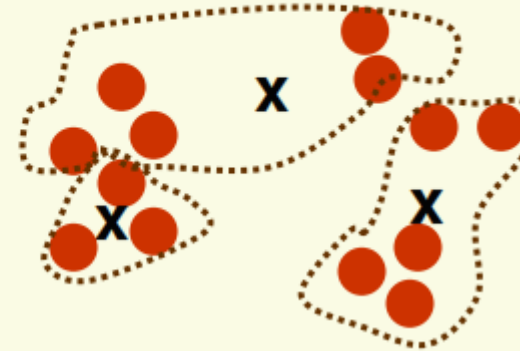
$k = 3$

1. Initialize

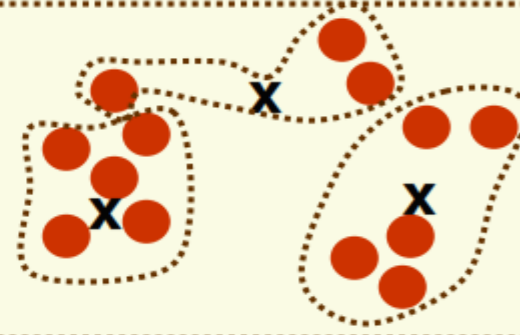
- pick k cluster centers arbitrary
- assign each example to closest center



2. compute sample means for each cluster



3. reassign all samples to the closest mean



4. if clusters changed at step 3, go to step 2