

## **More details of K-means clustering**

## K-means clustering algorithm

0. Start with initial guesses for cluster centers (centroids)
1. For each data point, find closest cluster center (partitioning step)
2. Replace each centroid by average of data points in its partition
3. Iterate 1+2 until convergence

(See Fig 14.4, 14.6)

Write  $x_i = (x_{i1}, \dots, x_{ip})$ :

If centroids are  $m_1, m_2, \dots, m_k$ , and partitions are

$c_1, c_2, \dots, c_k$ , then one can show that K-means converges to a *local* minimum of

$$\sum_{k=1}^K \sum_{i \in c_k} ||x_i - m_k||^2 \quad \text{Euclidean distance}$$

(within cluster sum of squares)

### **In practice:**

- Try many random starting centroids (observations) and choose solution with smallest of squares

### **How to choose K?**

- Difficult – details later

## Stepping back

- All clustering algorithms start with a dissimilarity measure for  $j^{th}$  feature

$d_j(x_{ij}, x_{i'j})$  and define

$$D(x_i, x_{i'}) = \sum_{j=1}^P d_j(x_{ij}, x_{i'j})$$

Usually  $d_j(x_{ij}, x_{i'j}) = (x_{ij} - x_{i'j})^2$

## Other possibilities:

- Correlation

$$\rho(x_i, x_{i'}) = \frac{\sum_j (x_{ij} - \bar{x}_i)(x_{i'j} - \bar{x}_{i'})}{\sqrt{\sum_j (x_{ij} - \bar{x}_i)^2 \sum_j (x_{i'j} - \bar{x}_{i'})^2}}$$

$\bar{x}_i$  = mean of observation  $i$

- If observations are standardized:

$$x_{ij} \leftarrow \frac{x_{ij} - \bar{x}_i}{\sqrt{\sum_j (x_{ij} - \bar{x}_i)^2}}$$

then  $2(1 - \rho(x_i, x_{i'})) = \sum_j (x_{ij} - x_{i,j'})^2$

So clustering via correlation  $\equiv$  clustering via Euclidean distance with standardized features

- **Categorical** features - usually coded as dummy variables

$$\begin{array}{lcl} \text{e.g. } X_1 = 1, 2 \text{ or } 3 & \rightarrow & (1 \ 0 \ 0) \\ & & (0 \ 1 \ 0) \\ & & \text{or } (0 \ 0 \ 1) \end{array}$$

- **Weighting** is also possible (see chapter 14)

## Partitioning (Clustering) Algorithms

- Group assignment function (“encoder”)  $C(i)$

$$C : 1, 2, \dots, N \rightarrow (1, 2, \dots, K)$$

- **Criterion:** choose  $C$  to minimize

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} d(x_i, x_{i'})$$

(within cluster scatter)

**Fact:**

- $K$ -means minimizes  $W(C)$  when  $D = \|x_i - x_{i'}\|^2$

$$\begin{aligned} W(C) &= \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} \|x_i - x_{i'}\|^2 \\ &= \sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - \bar{x}_k\|^2 \end{aligned}$$

- $K$ -means solves *enlarged* problem:

$$\min_{C, m_1 \dots m_k} \sum_k \sum_{C(i)=k} \|x_i - m_k\|^2$$

to find assignment function  $C$