

Clustering using the kmeans algorithm

Exercise 1 Partition and matrix

Consider the Iris data set. Write a R code wich produces the partition matrix. Compute the gravity centers of the quantitative variables in the three classes using a matrix formula.

Exercise 2 The bell number

1. Show that the number of partition of n objects verifies

$$B_{n+1} = \sum_{k=0}^n C_k^n B_k$$

2. Compute manually the bell number for 1,2,3,4,5,6 objects.
3. Write a R program which computes the Bell number for n objects.

Exercise 3 Between-Within Variance relation

Consider n points from \mathbb{R}^p with a partition into K classes of size n_1, \dots, n_k . Let us note $\hat{\mu}_k$ the gravity center of class k and $\hat{\mu}$ the gravity center of the entire cloud of points. Show that

$$\sum_k \sum_{i \in k} \|\mathbf{x}_i - \hat{\mu}_k\|^2 + \sum_k n_k \|\hat{\mu}_k - \hat{\mu}\|^2 = \sum_i \|\mathbf{x}_i - \hat{\mu}\|^2$$

Exercise 4 Clustering of the crabs (library MASS)

1. Load the `crabs` dataset form library `MASS`.
2. Plot the dateset using `pairs()` with a color for each specy and a different symbol per sex.
3. Cluster the dataset reduced to its quantitative variables into four cluster using the `kmeans`.
4. Run the algorithm with 1000 different initialization and keep track of the within sum of squares.
5. Comment the result.
6. Divide all quantitative variable by the most correlated variable to produce a new dataset.
7. Compare the partitions obtained using the `kmeans` with the 'natural' partition. Comment.
8. Try to cluster the data in 1 to 20 groups. Plot the within sum of squares in function of the number of clusters. Comment the figure.