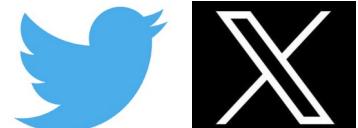
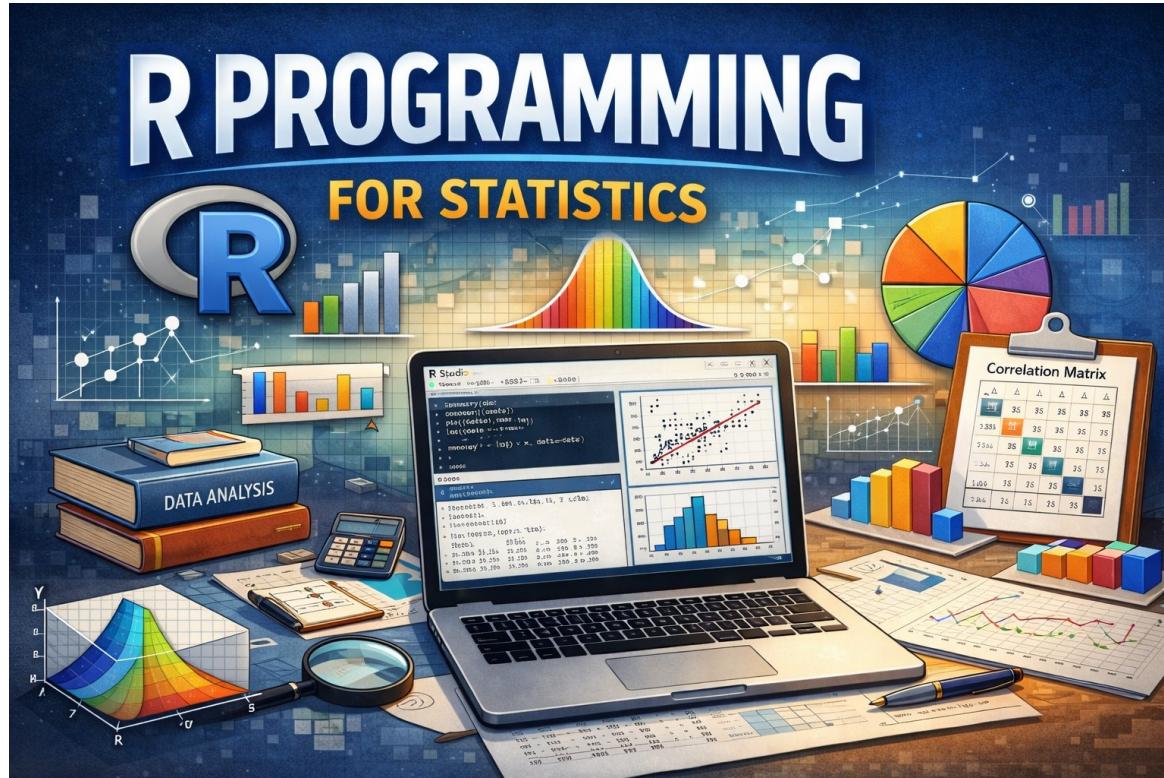




Latvian Institute of
Organic Synthesis

Basic Clustering in R

Nikolay Oskolkov, MRG Group Leader, LIOS, Riga, Latvia
R course, 12.02.2026



@NikolayOskolkov



@oskolkov.bsky.social



Personal homepage:
<https://nikolay-oskolkov.com>

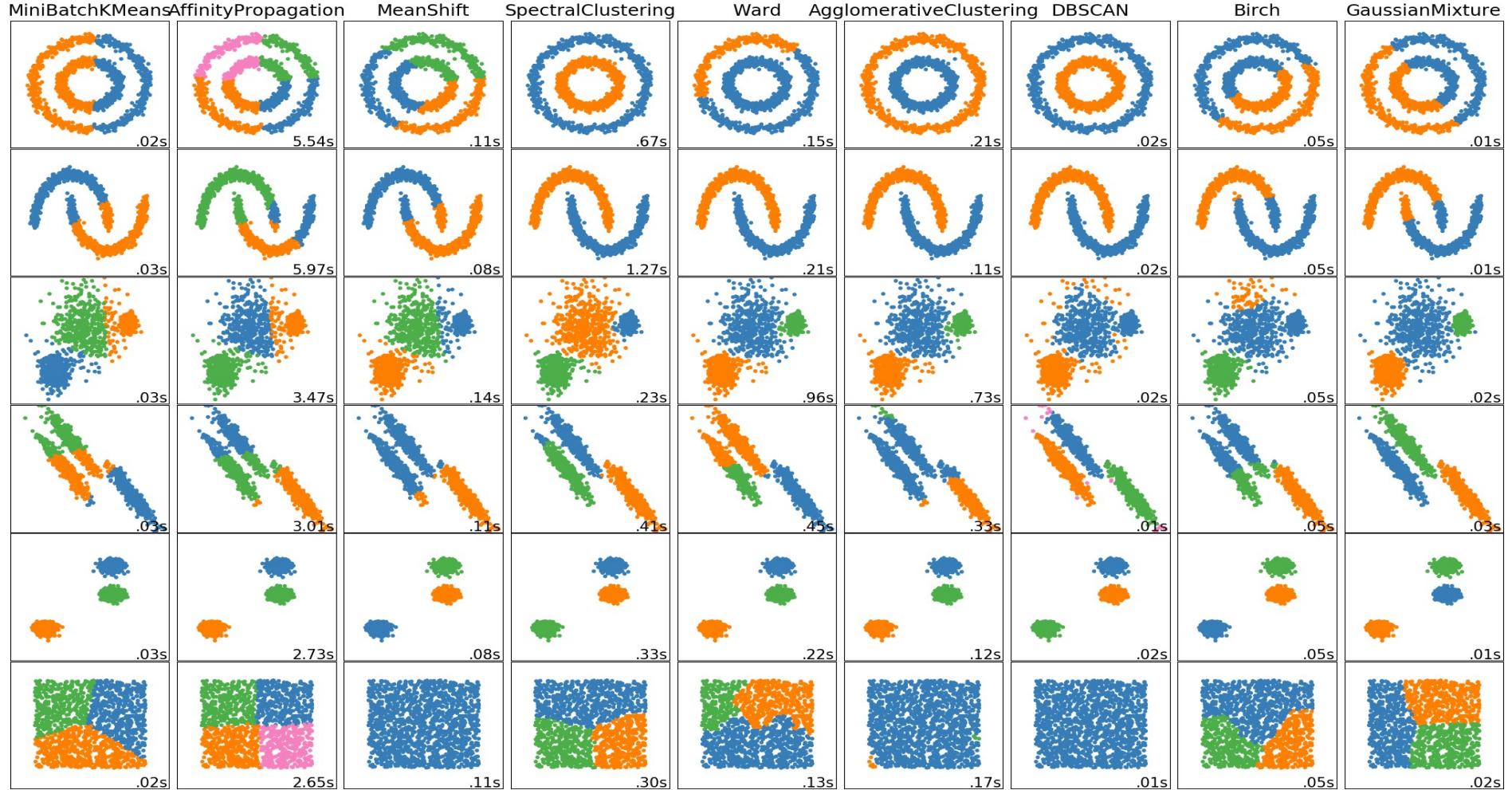


Session content

Topics we'll cover in this session:

- 1) Hierarchical agglomerative clustering
- 2) Partitioning (divisive) clustering
- 3) K-means and Gaussian Mixture Model (GMM)
- 4) Coding K-means and GMM from scratch in R and Python

Huge variety of clustering methods available

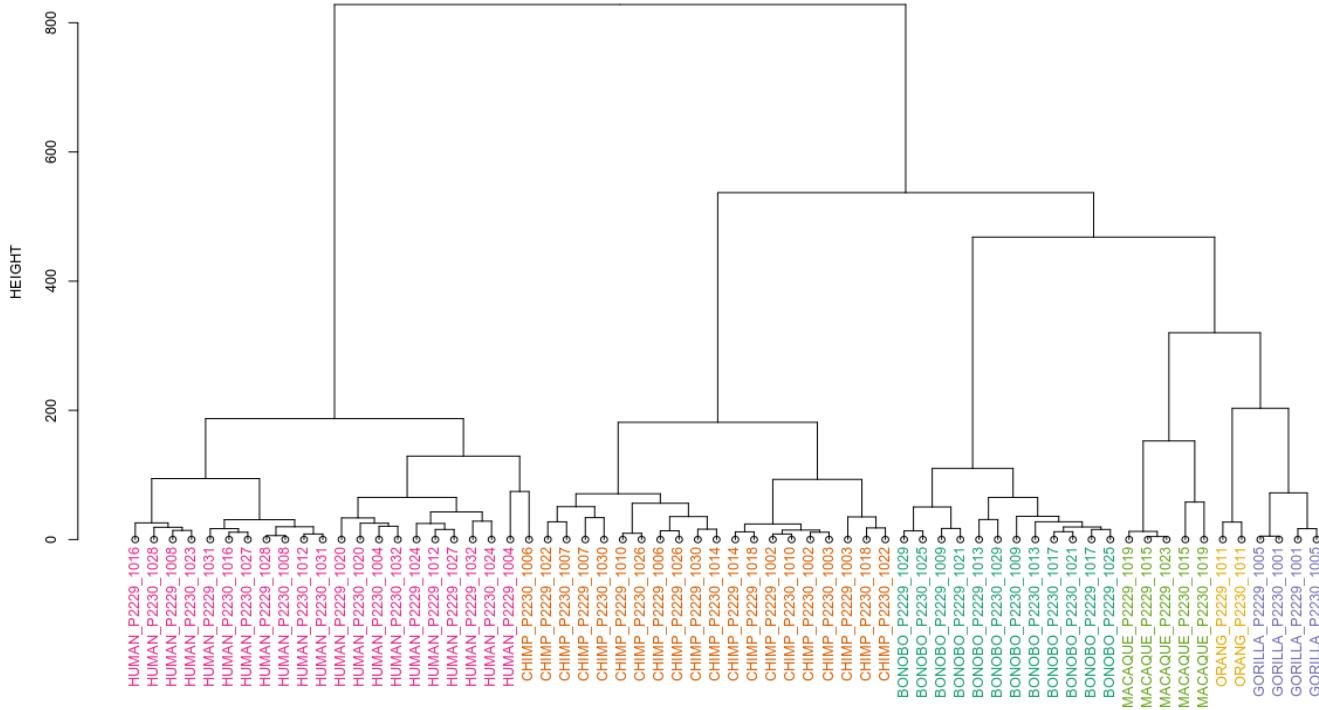


1) Hierarchical Clustering, 2) Partitioning Clustering, 3) Density-Based Clustering, 4) Graph-Based Clustering

Hierarchical clustering

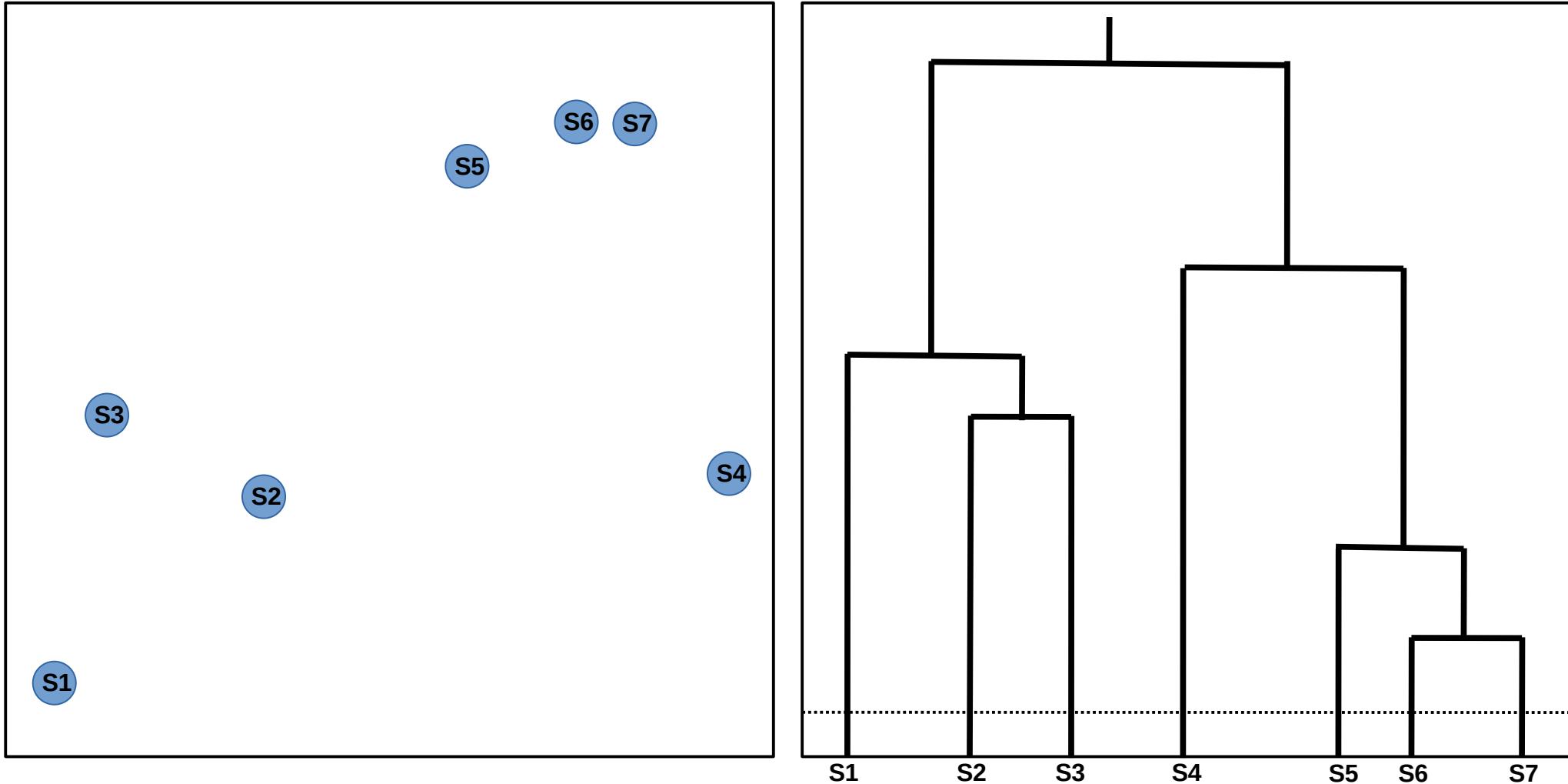
What is hierarchical clustering?

Clusters are nested (higher-ranked clusters contain lower-ranked clusters),
in other words tree-like structure (dendrogram)

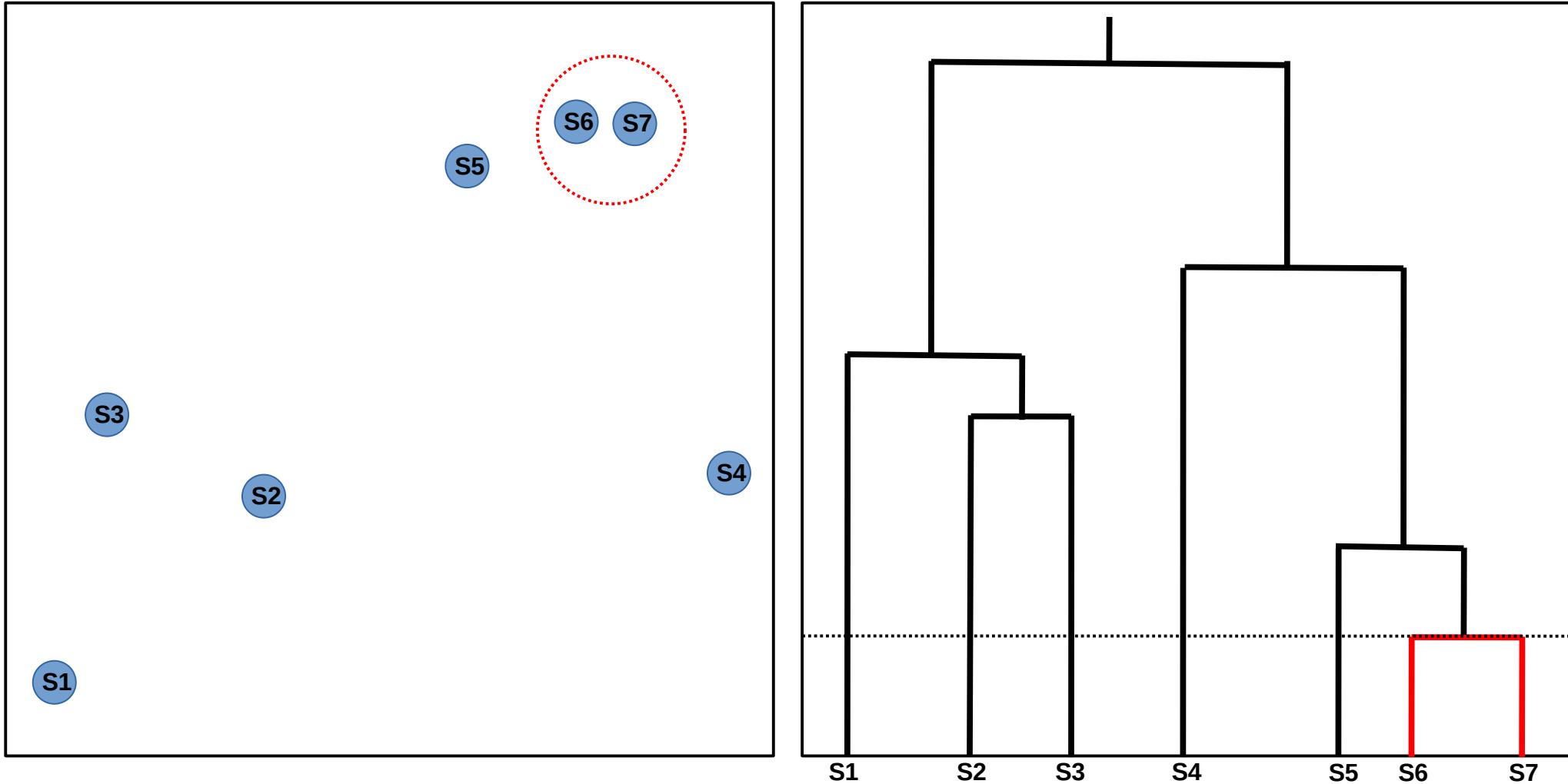


There are agglomerative (bottom-up approach) and divisive (top-down approach) hierarchical types of hierarchical clustering

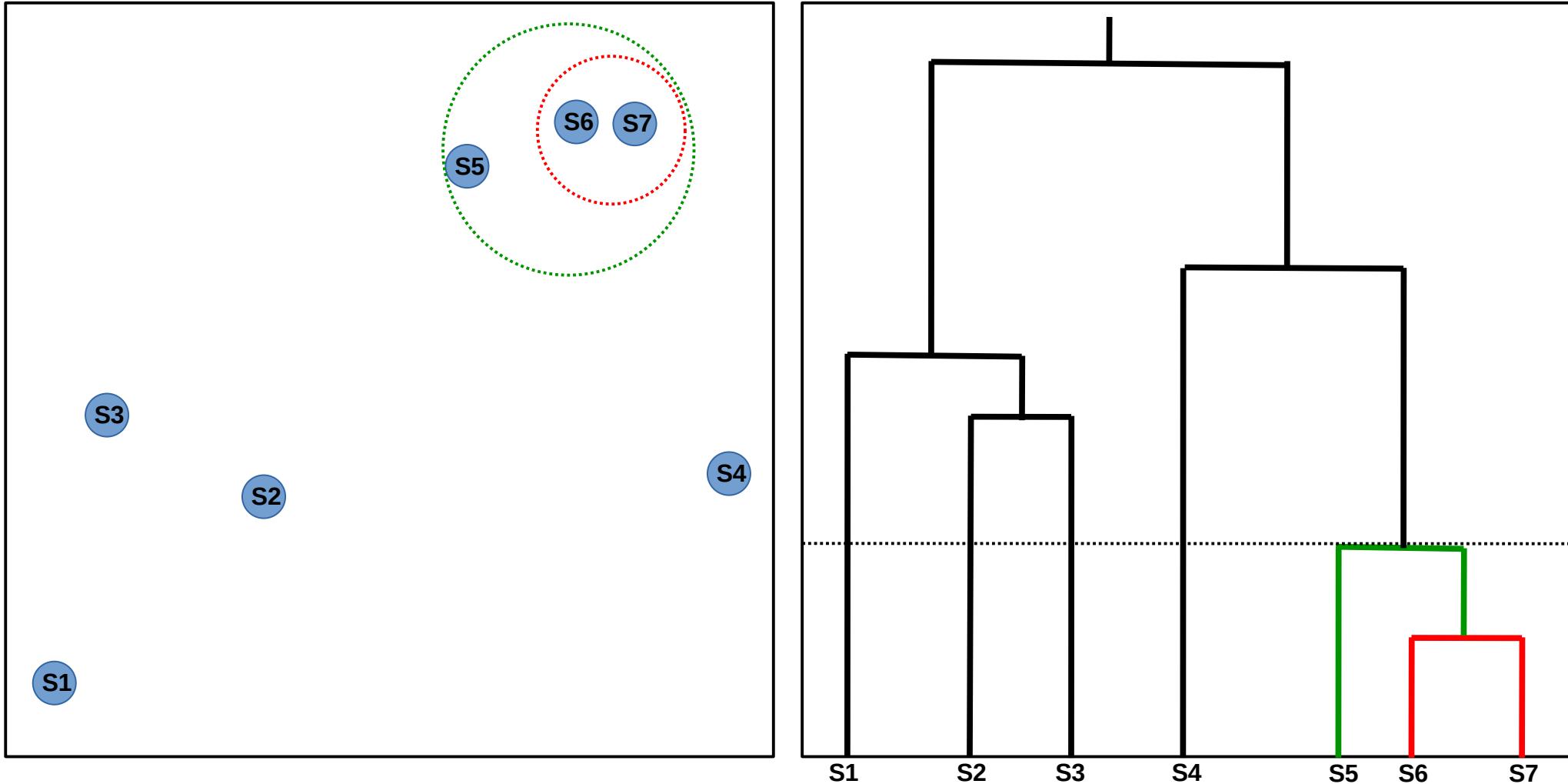
Agglomerative hierarchical clustering



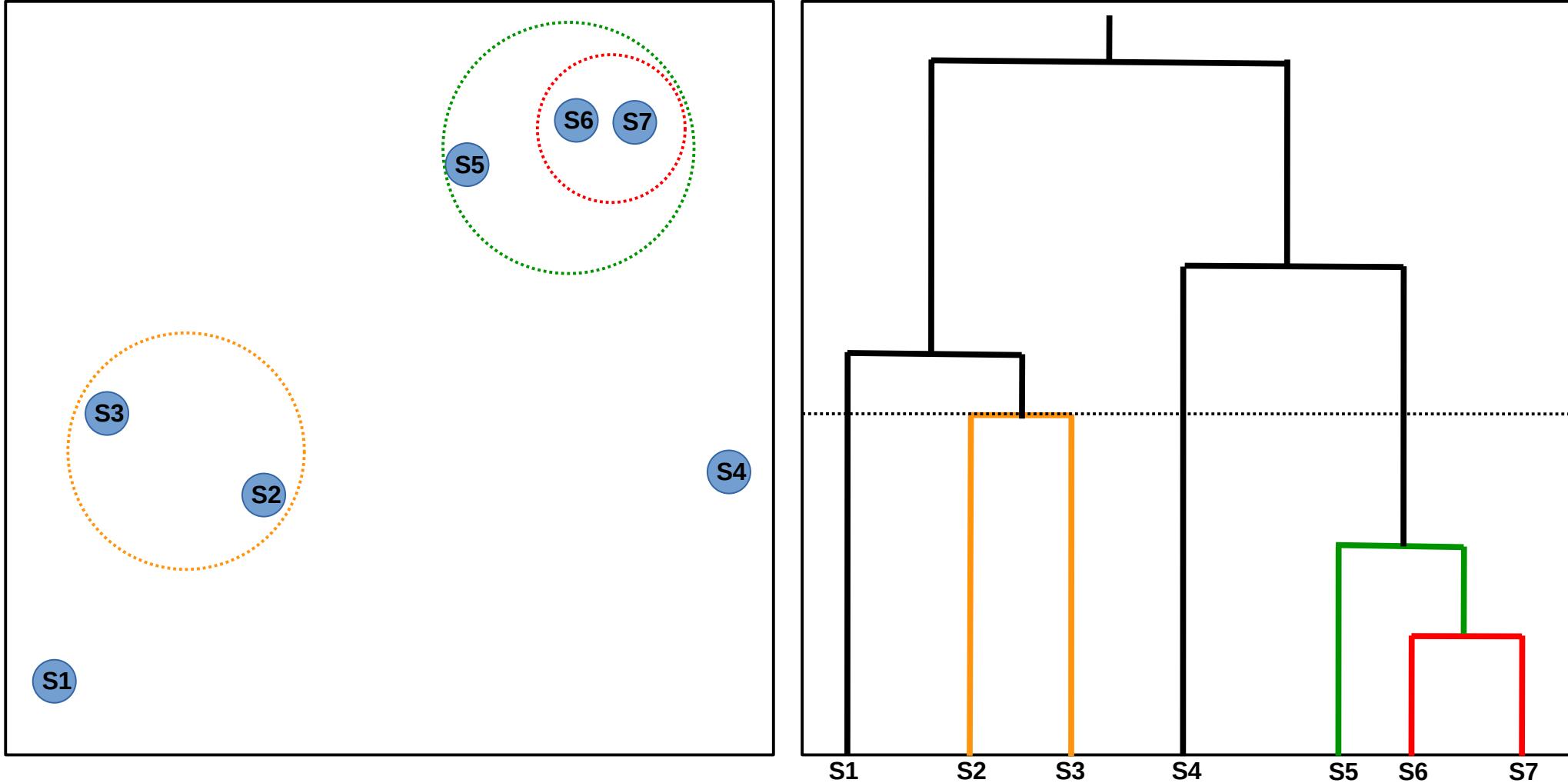
Agglomerative hierarchical clustering



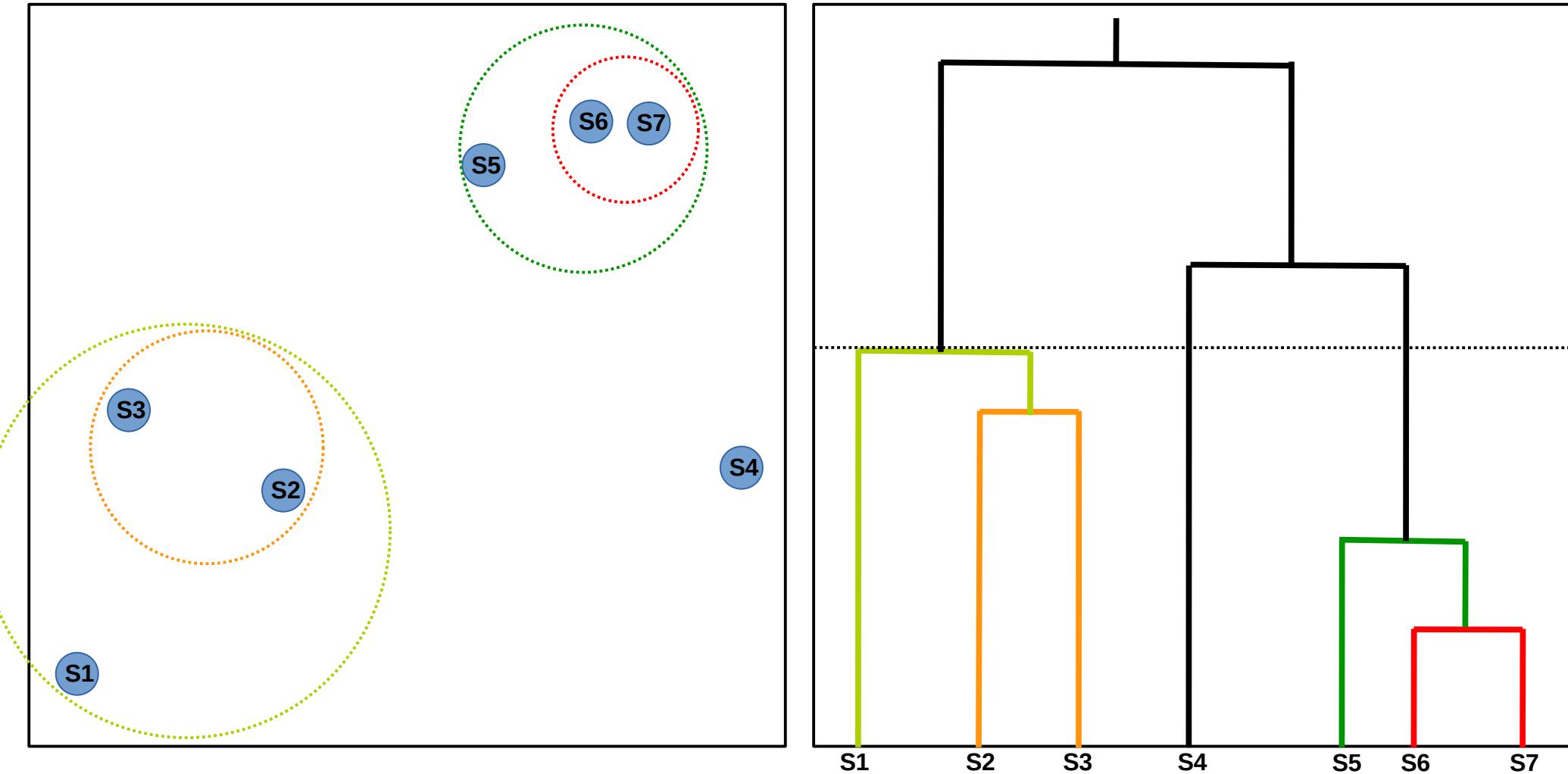
Agglomerative hierarchical clustering

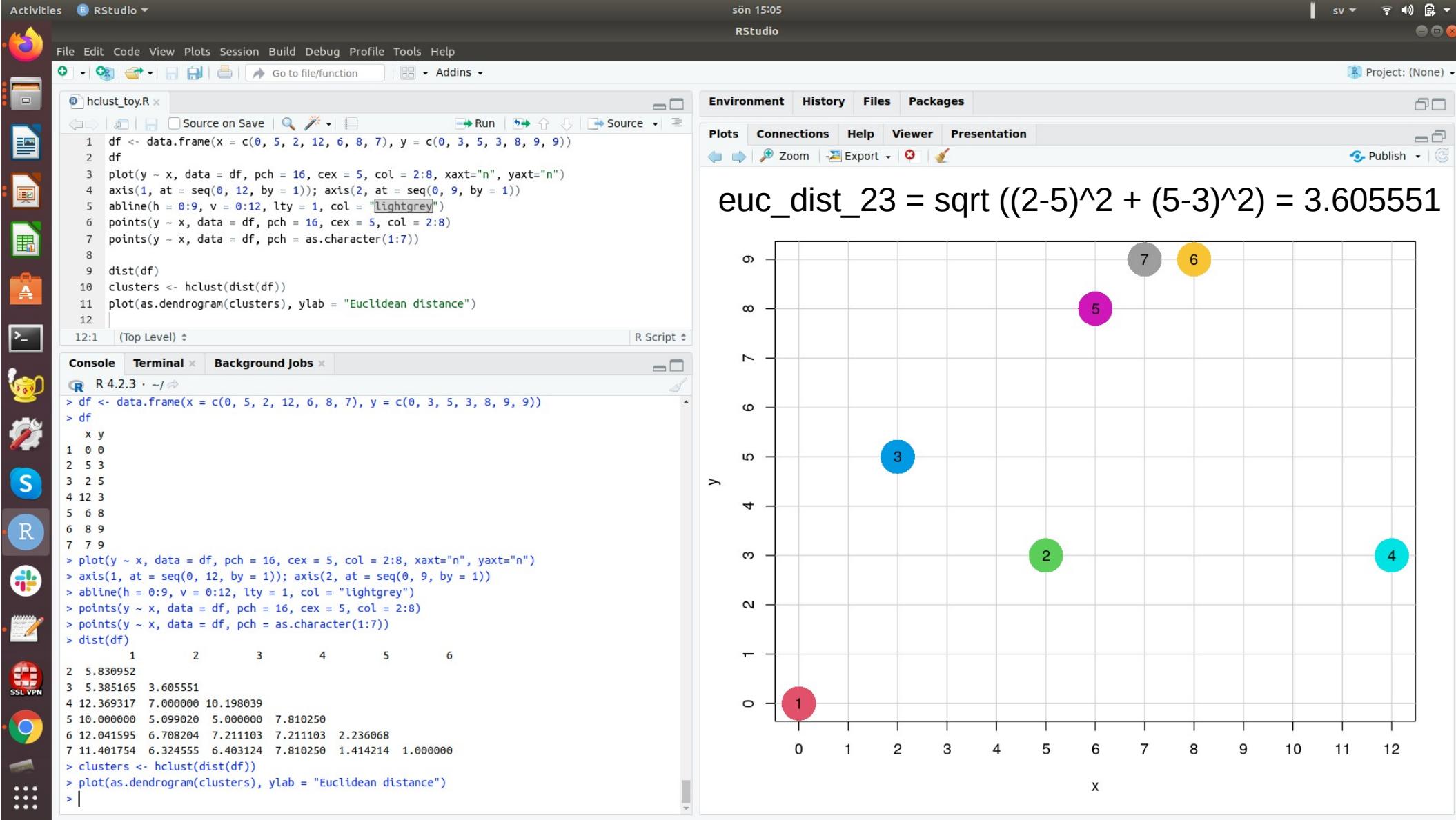


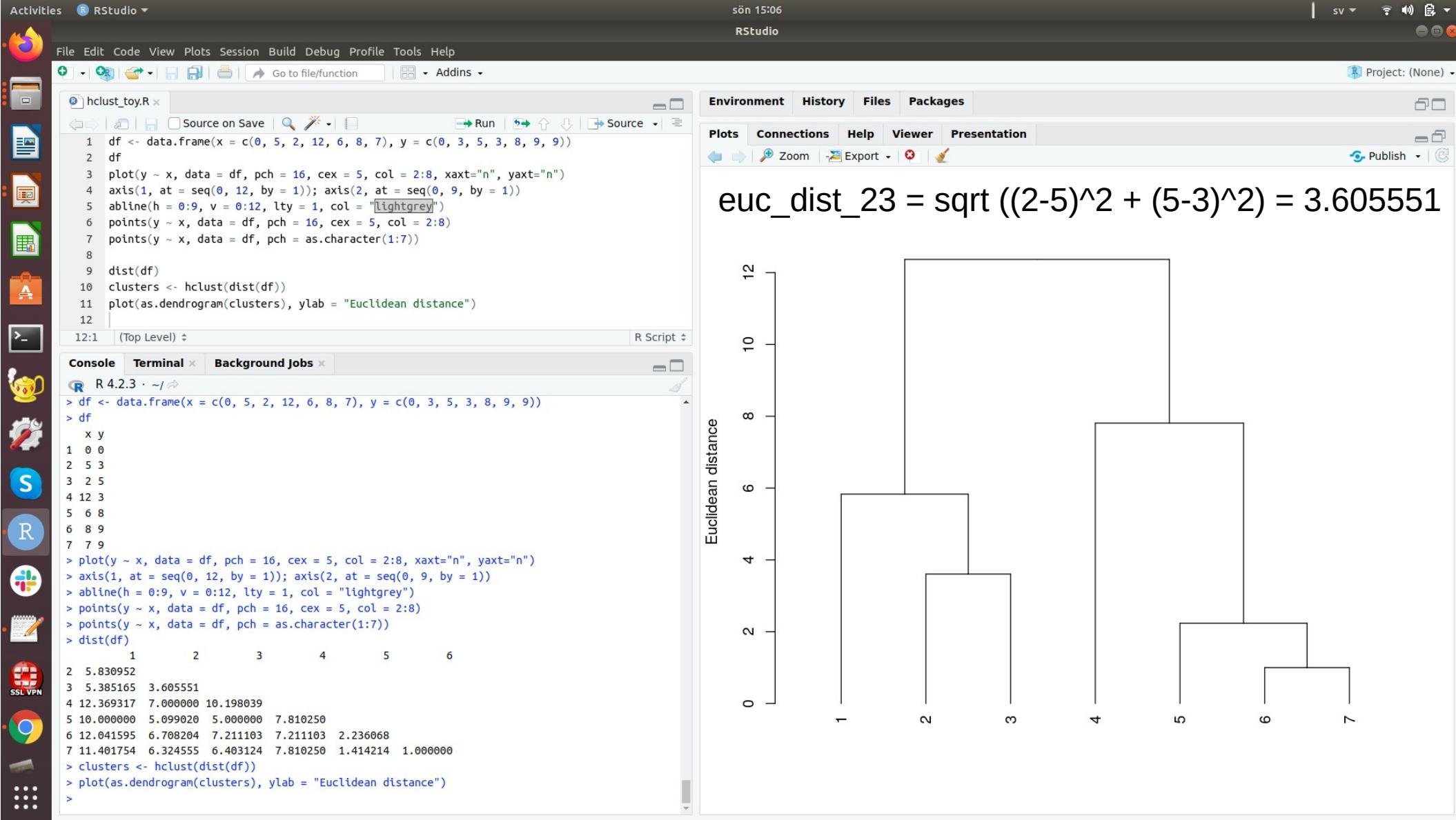
Agglomerative hierarchical clustering



Agglomerative hierarchical clustering







Types of hierarchical clustering

Add sample to cluster by:

- 1) Closest neighbor in cluster (single linkage)
- 2) Farthest neighbor in cluster (complete linkage)
- 3) Median / mean of cluster (UPGMA / UPGMC)
- 4) Least change in variation (Ward)

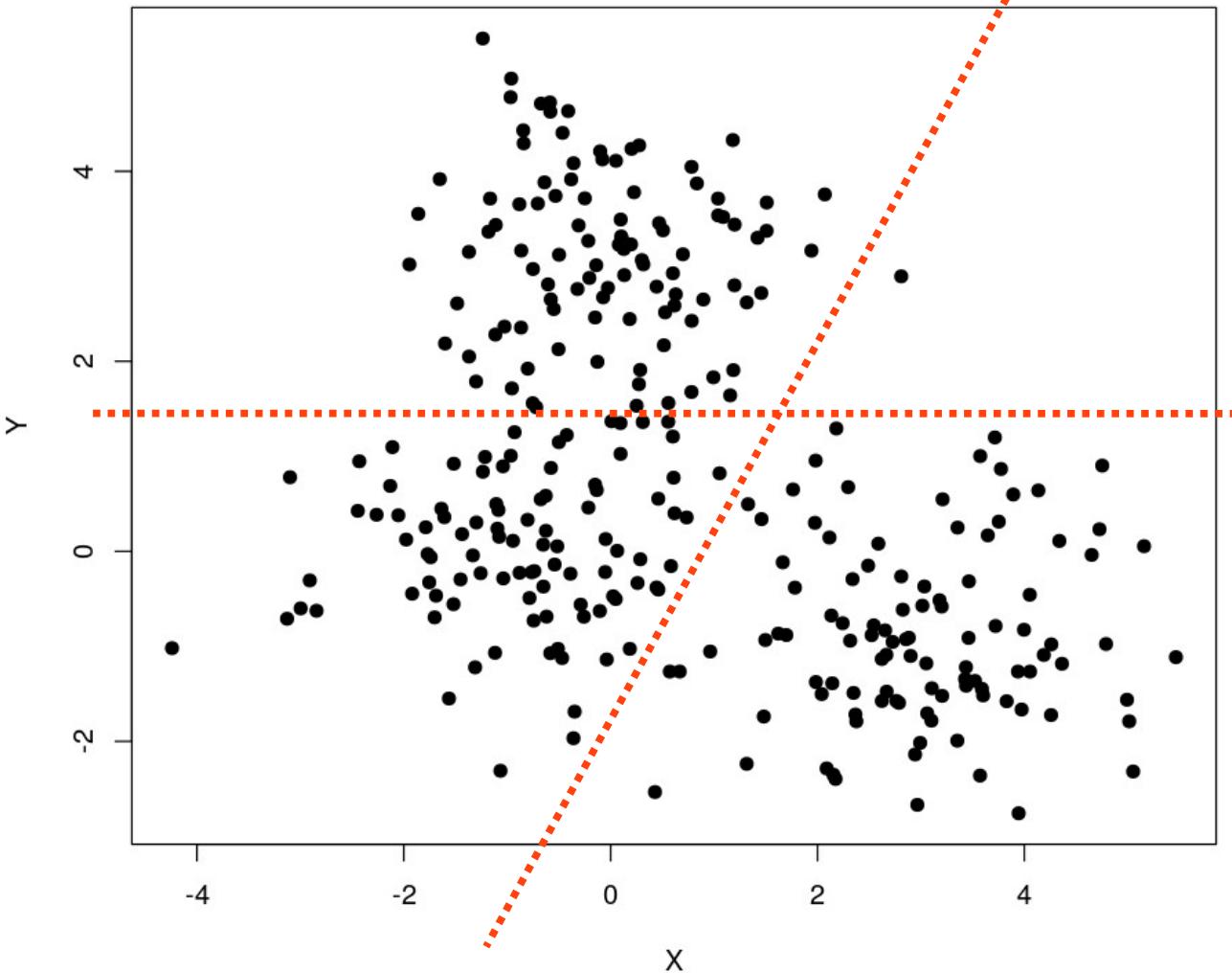
Disadvantages of hierarchical clustering (no matter type):

Very sensitive to noise in data



Partitioning clustering

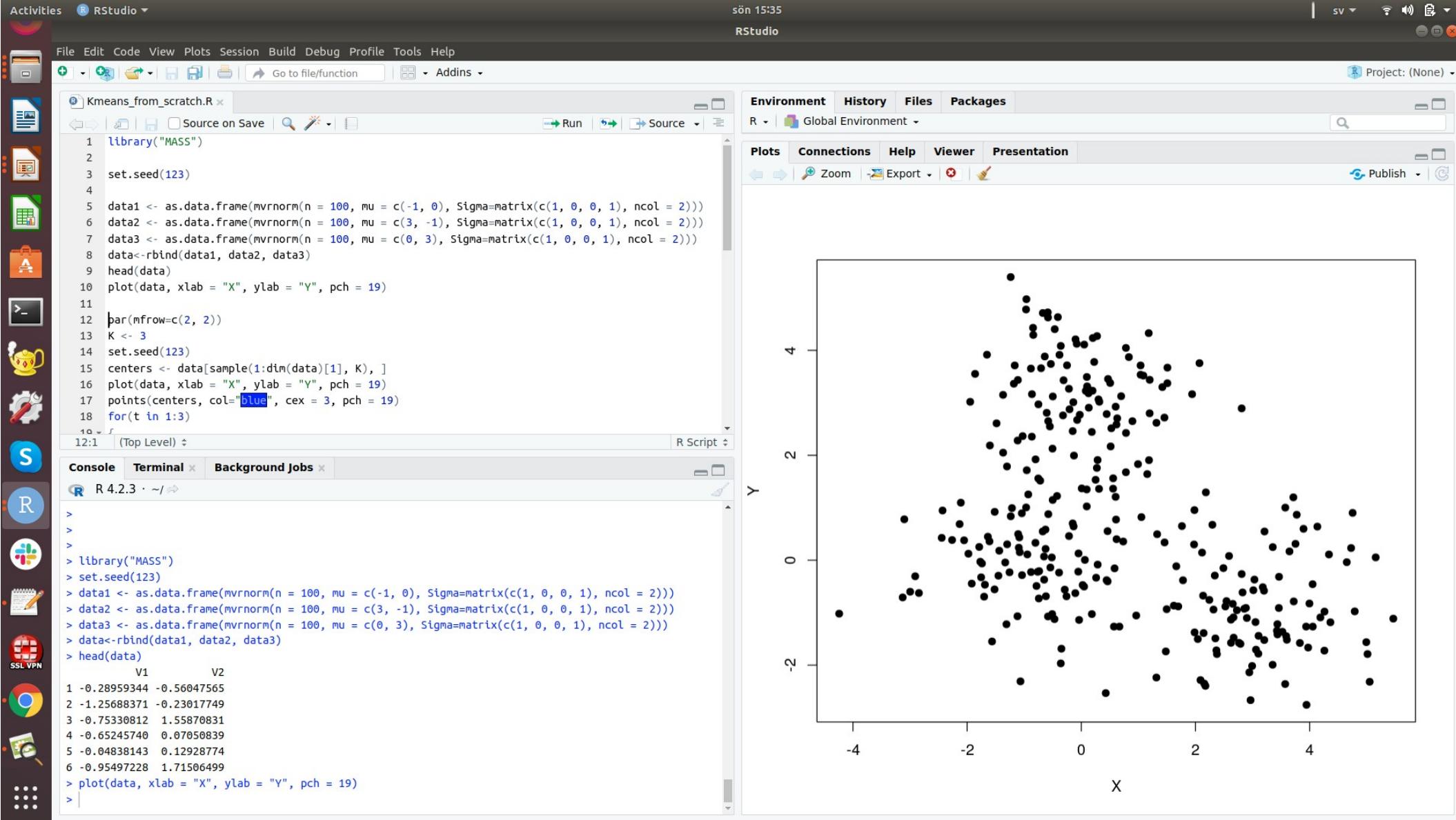
Partitioning principle



Assign points to
 $K = 2, 3, 4, 5, \dots$
groups (clusters)

No dendrogram!

K-means, C-means,
PAM, GMM, ...



Activities RStudio ▾ | sv ▾ | WiFi | Sound | Close

sön 15:36
RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Kmeans_from_scratch.R

```
12 par(mfrow=c(2, 2))
13 K <- 3
14 set.seed(123)
15 centers <- data[sample(1:dim(data)[1], K), ]
16 plot(data, xlab = "X", ylab = "Y", pch = 19)
17 points(centers, col="blue", cex = 3, pch = 19)
18 for(t in 1:3)
19 {
20 my_labels <- vector()
21 for(i in 1:dim(data)[1])
22 {
23 my_dist <- vector()
24 for(j in 1:K)
25 {
26 my_dist[j] <- sqrt((data[i, 1] - centers[j, 1])^2 + (data[i, 2] - centers[j, 2])^2)
27 }
28 my_labels[i] <- which.min(my_dist)
29 }
30 plot(data, xlab = "X", ylab = "Y", pch = 19)
31 points(data, col = my_labels, pch = 19)
32 points(centers, col = "blue", cex = 3, pch = 19)
33 |
34 new_centers <- list()
35 for(i in unique(my_labels))
36 {
37 new_centers[[i]] <- colMeans(data[my_labels == i,])
38 }
39 centers <- Reduce("rbind", new_centers)
40 }
```

33:1 (Top Level) ▾ R Script

Console Terminal × Background Jobs ×

```
R 4.2.3 · ~/🔗
+ new_centers <- list()
+ for(i in unique(my_labels))
+ {
+ new_centers[[i]] <- colMeans(data[my_labels == i,])
+ }
+ centers <- Reduce("rbind", new_centers)
+ }
> |
```

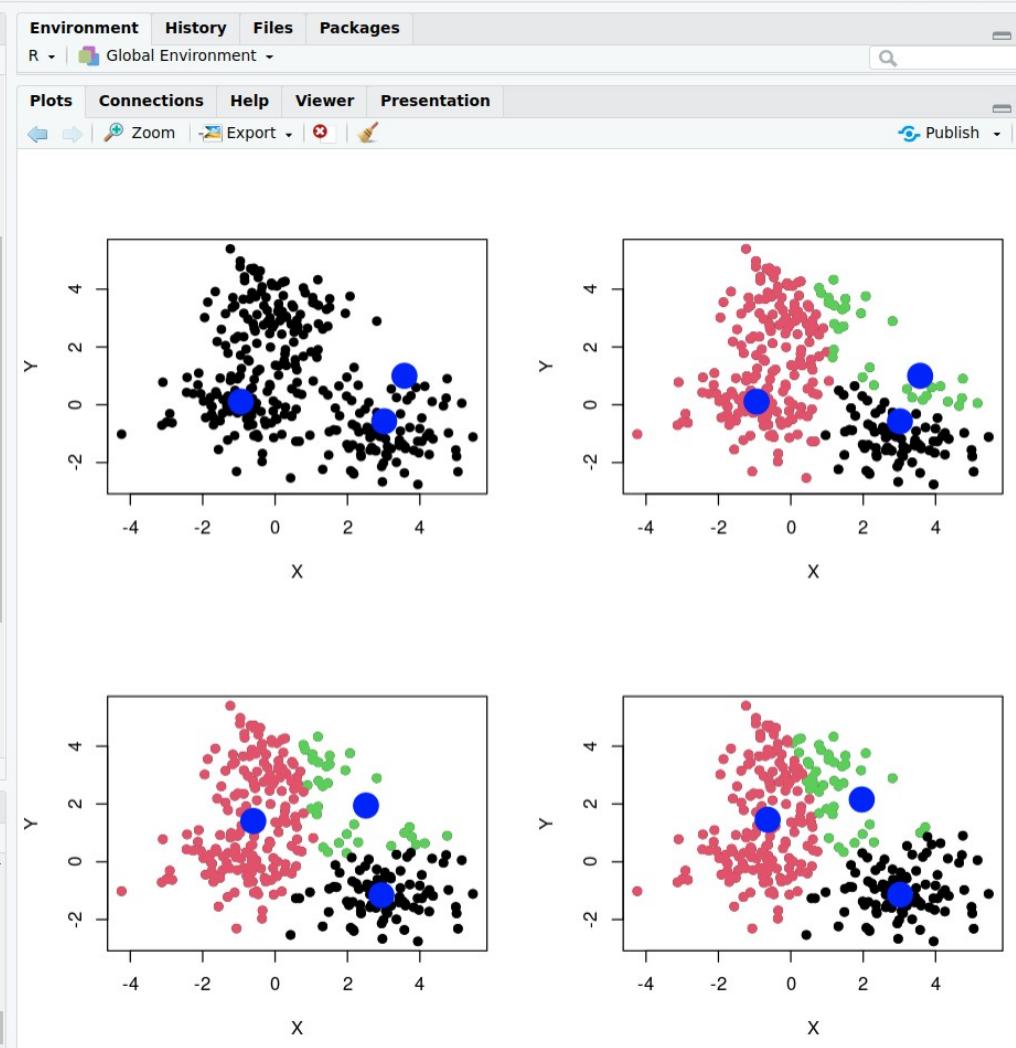
Environment History Files Packages

R Global Environment

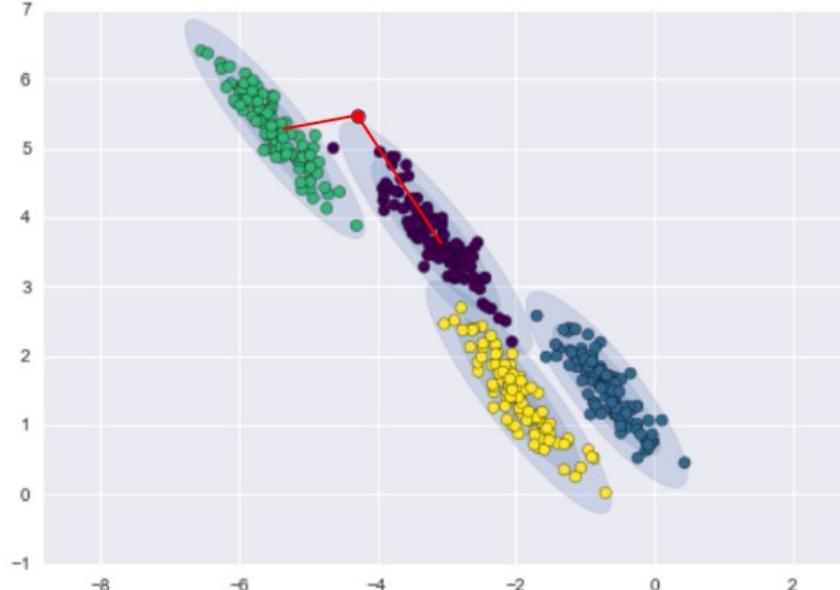
Plots Connections Help Viewer Presentation

Zoom Export Publish

Four scatter plots showing the k-means clustering process. The first plot shows the initial data points and three blue centroids. Subsequent plots show the iterative assignment of data points to clusters (red and green) and the subsequent shift of the centroids (blue circles) towards the cluster means.



Gaussian Mixture Model (GMM)



GMM is a probabilistic clustering

Expectation-Maximization (EM) algorithm:

$$L(X|\mu_k, \sigma_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(X - \mu_k)^2}{2\sigma_k^2}}$$

Once we computed the likelihood $L(X|\mu, \sigma)$, which we assume follows the Normal distribution, of observing the data given the parameters (means and variances), we can compute the posterior probability of data point assignment to each cluster using the Bayes rule.

$$p_k(\mu_k, \sigma_k|X) = \frac{L(X|\mu_k, \sigma_k)\phi_k(\mu_k, \sigma_k)}{\sum_k L(X|\mu_k, \sigma_k)\phi_k(\mu_k, \sigma_k)}$$

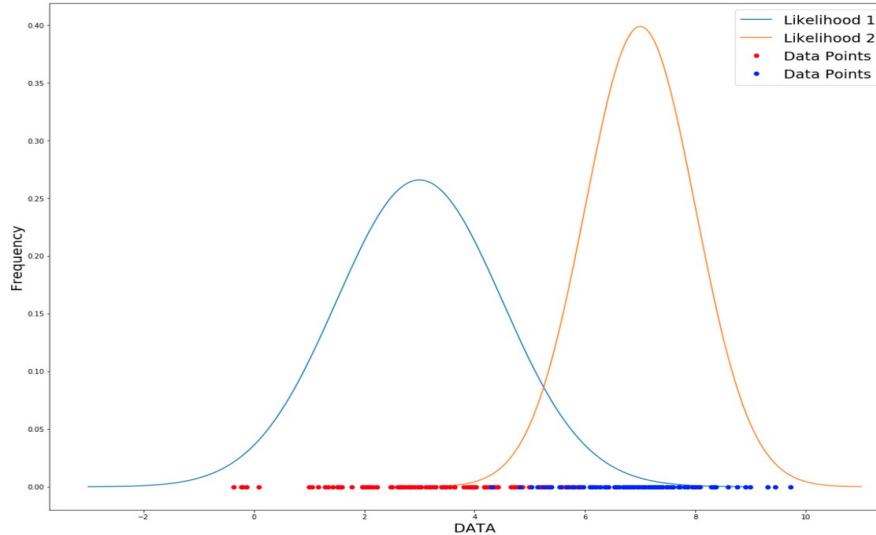
Next, differentiating the posterior probability with respect to the parameters, means and variances, i.e. by maximizing the posterior probability of the parameters given the data, we can derive a rule for updating the means and variances of the clusters.

$$\mu_k = \frac{\sum_i X_i p_k(\mu_k, \sigma_k|X_i)}{\sum_i p_k(\mu_k, \sigma_k|X_i)}, \quad \sigma_k^2 = \frac{\sum_i (X_i - \mu_k)^2 p_k(\mu_k, \sigma_k|X_i)}{\sum_i p_k(\mu_k, \sigma_k|X_i)}, \quad \phi_k = \frac{1}{N} \sum_i p_k(\mu_k, \sigma_k|X_i)$$

Disadvantages of partitioning clustering (no matter type):

- 1) Assumes spherical / ellipsoidal symmetry of clusters
- 2) Number of clusters has to be defined a-priori

GMM from scratch in Python



```

import random
import numpy as np
from scipy.stats import norm
import matplotlib.pyplot as plt

n = 100
mu1 = 3; sigma1 = 1.5
mu2 = 7; sigma2 = 1

x1 = np.random.normal(mu1, sigma1, n)
x2 = np.random.normal(mu2, sigma2, n)
X = np.array(list(x1) + list(x2))
np.random.shuffle(X)
print("Dataset shape:", X.shape)

x1_approx = np.linspace(mu1 - 4*sigma1, mu1 + 4*sigma1, 10000)
x2_approx = np.linspace(mu2 - 4*sigma2, mu2 + 4*sigma2, 10000)
X_approx = np.array(list(x1_approx) + list(x2_approx))

def pdf(data, mean, sigma):
    return (1 / np.sqrt(2 * np.pi * sigma**2)) * np.exp(-(np.square(data - mean) / (2 * sigma**2)))

plt.figure(figsize = (20,15))
# plt.plot(x1_approx, norm.pdf(x1_approx, mu1, sigma1))
# plt.plot(x2_approx, norm.pdf(x2_approx, mu2, sigma2))
plt.plot(x1_approx, pdf(x1_approx, mu1, sigma1))
plt.plot(x2_approx, pdf(x2_approx, mu2, sigma2))

plt.plot(x1, list(0 for i in range(0, len(list(x1))), 'ro'))
plt.plot(x2, list(0 for i in range(0, len(list(x2))), 'bo'))

plt.ylabel('Frequency', fontsize = 20)
plt.xlabel('DATA', fontsize = 20)
plt.legend(['Likelihood 1', 'Likelihood 2', 'Data Points 1', 'Data Points 2'], loc = 'upper right', fontsize = 20)
plt.show()

```

Coding Expectation-Maximization (EM) algorithm from scratch in Python:

```

k = 2
np.random.seed(123)
means = np.random.choice(X, k)
sigmas = np.random.random_sample(size = k)
priors = np.ones((k)) / k
print(means, sigmas, priors)

from sklearn import cluster, datasets, mixture
from scipy.stats import multivariate_normal

N_iter = 100
bins = np.linspace(np.min(X_approx), np.max(X_approx), 10000)
for step in range(N_iter + 1):
    if step % 10 == 0:

        plt.figure(figsize = (20, 15))
        plt.ylabel('Frequency', fontsize = 20)
        plt.xlabel('DATA', fontsize = 20)
        plt.title("Iteration {}".format(step), fontsize = 20)

        plt.plot(x1_approx, pdf(x1_approx, mu1, sigma1))
        plt.plot(x2_approx, pdf(x2_approx, mu2, sigma2))

        plt.plot(x1, list(0 for i in range(0, len(list(x1))), 'ro'))
        plt.plot(x2, list(0 for i in range(0, len(list(x2))), 'bo'))

        print(means[0], means[1])
        print(sigmas[0], sigmas[1])

        plt.plot(bins, pdf(bins, means[0], sigmas[0]), color = 'darkgreen')
        plt.plot(bins, pdf(bins, means[1], sigmas[1]), color = 'darkgreen')

        plt.legend(['Likelihood 1', 'Likelihood 2', 'Data Points 1', 'Data Points 2', 'Approximation'],
                  loc = 'upper right', fontsize = 20)
        plt.show()

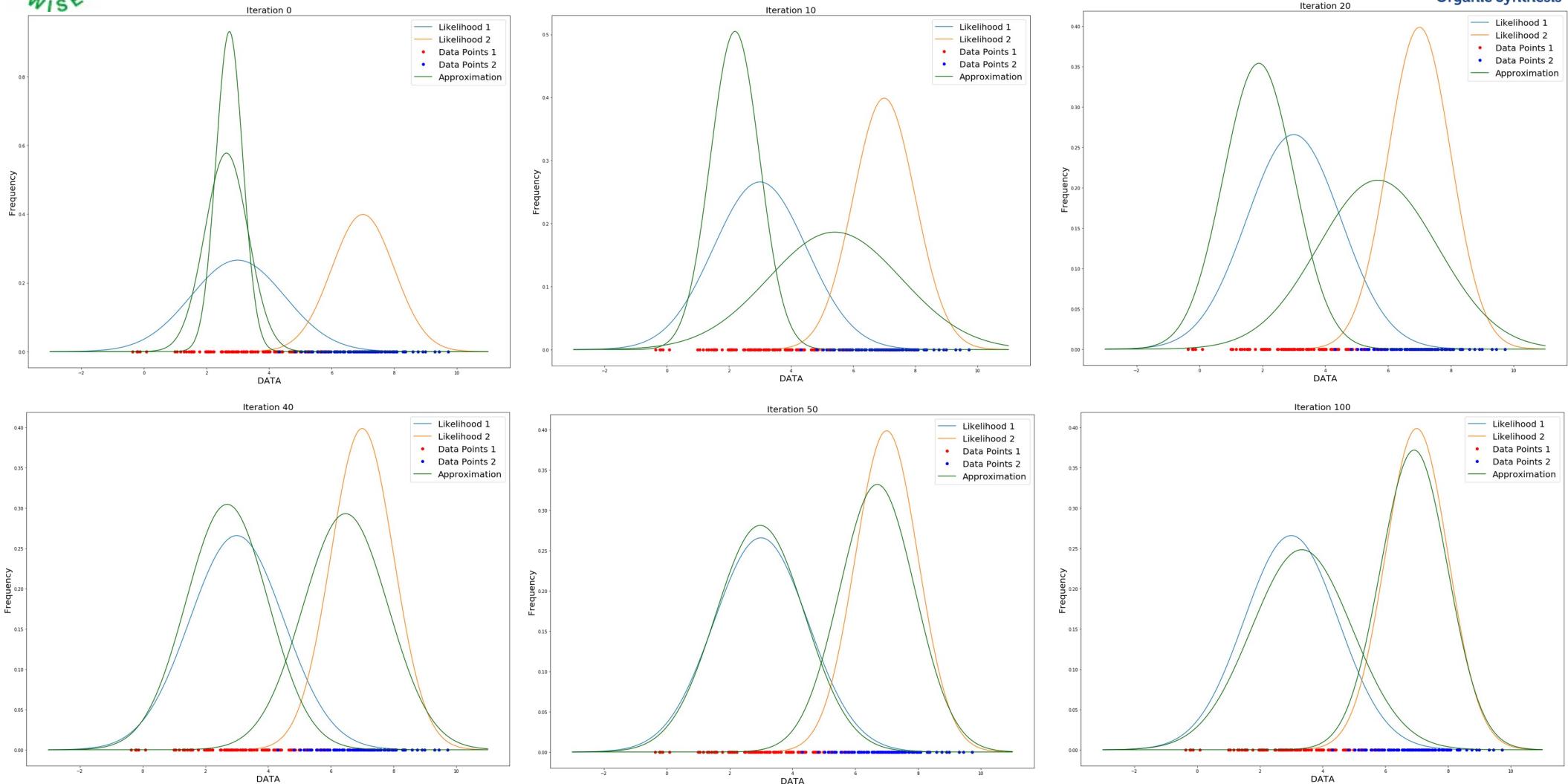
    # Expectation step
    likelihood = []
    for j in range(k):
        likelihood.append(pdf(X, means[j], sigmas[j]))

    # Maximization step
    posterior = []
    for j in range(k):
        posterior.append((likelihood[j] * priors[j]) / (np.sum([likelihood[i] * priors[i] for i in range(k)]),
                                                       axis = 0))

        means[j] = np.sum(posterior[j] * X) / np.sum(posterior[j])
        sigmas[j] = np.sqrt(np.sum(posterior[j] * np.square(X - means[j])) / np.sum(posterior[j]))
        priors[j] = np.mean(priors[j])

```

GMM: iteration by iteration



Session summary

Take home messages of the session:

- 1) Hierarchical clustering involves merging pairs of most similar points and building a dendrogram of their relationships, it is sensitive to noise in the data
- 2) Partitioning clustering implies dividing manifold of data points into groups of most similar points
- 3) K-means and Gaussian Mixture Model (GMM) assume spherical / ellipsoidal symmetry of clusters and demand the number of clusters to be known a-priori



Acknowledgments: LIOS + TARGETWISE



2027
National
Development Plan



Funded by
the European Union