

BIOINFORMATICS

How to cluster biological data?

Marco Beccuti

Università degli Studi di Torino

Dipartimento di Informatica

May 2019



Outline

- 1 Gene expression and clustering;
- 2 Clustering as optimization problem;
- 3 k-means clustering;
- 4 Hierarchical Clustering.

Chapter 8 in ***Bioinformatics Algorithms: An active Learning Approach (Vol.2)***.



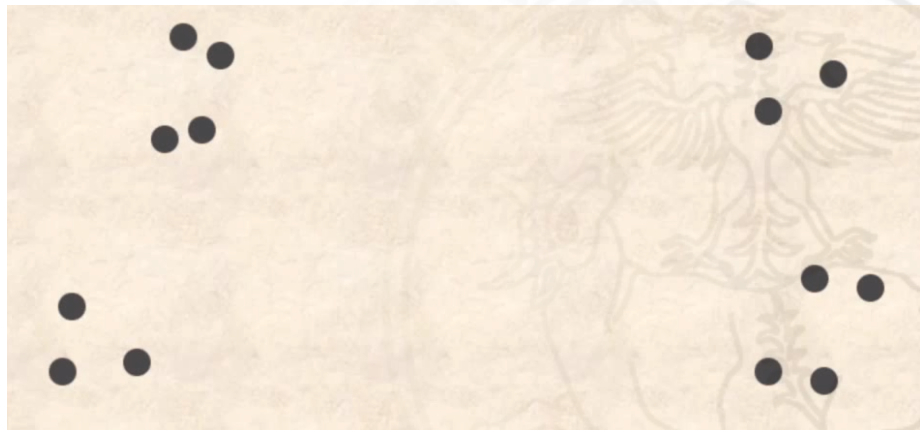
Part 1

Hierarchical Clustering

Hierarchical Clustering

Stratification of Clusters

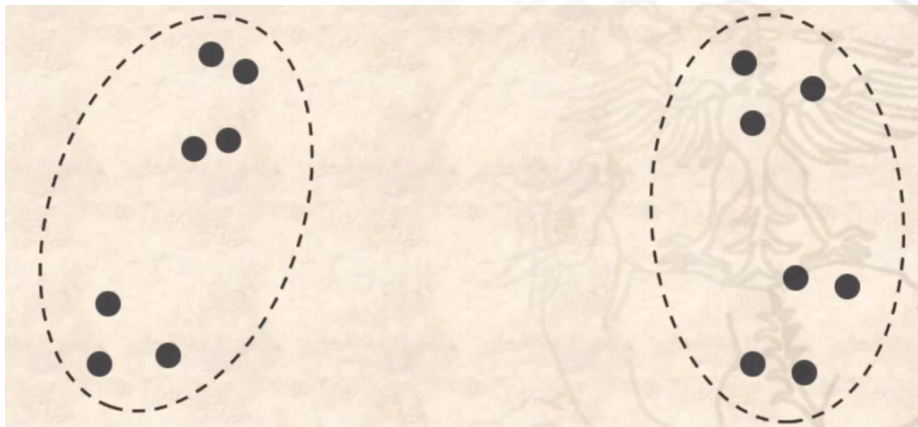
- Clusters often have **sub-cluster**, which have sub-clusters, and so on.



Hierarchical Clustering

Stratification of Clusters

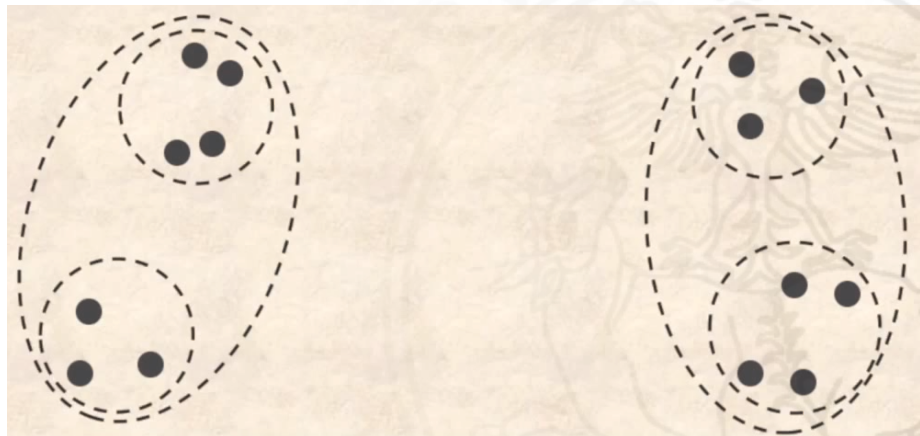
- Clusters often have **sub-cluster**, which have sub-clusters, and so on.



Hierarchical Clustering

Stratification of Clusters

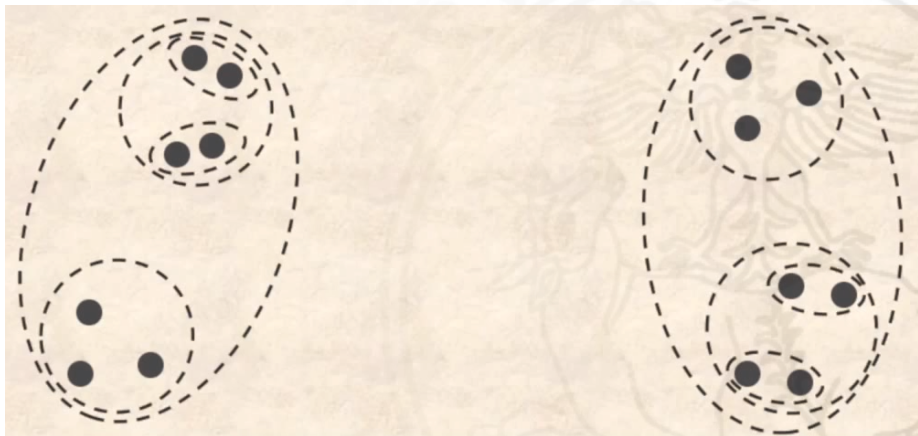
- Clusters often have **sub-cluster**, which have sub-clusters, and so on.



Hierarchical Clustering

Stratification of Clusters

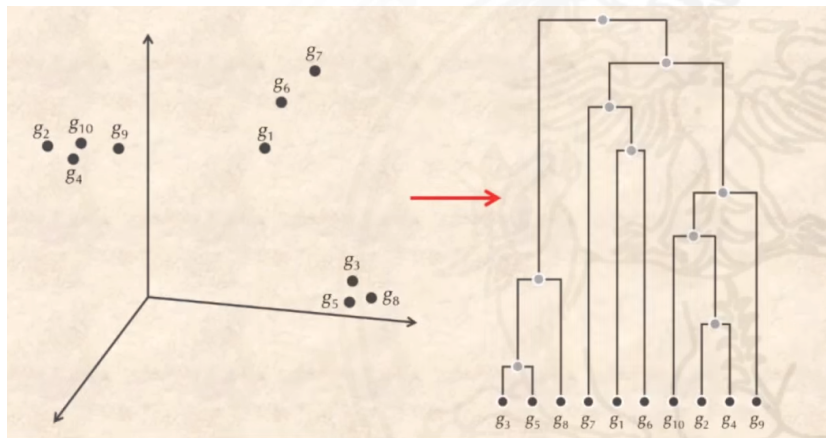
- Clusters often have **sub-cluster**, which have sub-clusters, and so on.



Hierarchical Clustering

From Data to a Tree

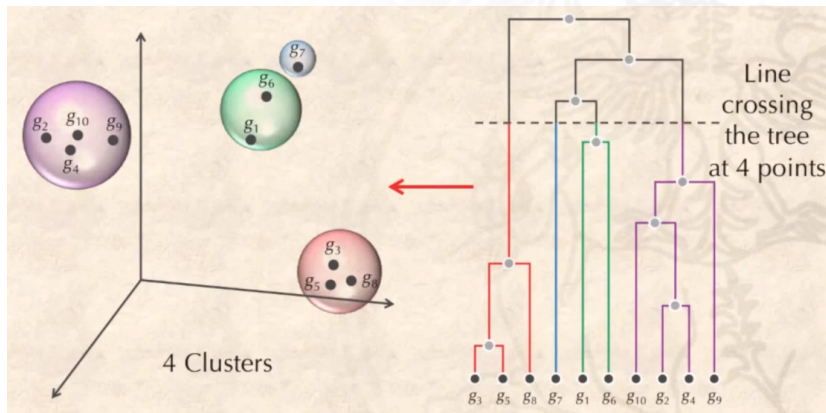
- To capture stratification, the **hierarchical clustering** algorithm organizes n data points into a tree (namely **Dendrogram**).



Hierarchical Clustering

From Data to a Tree

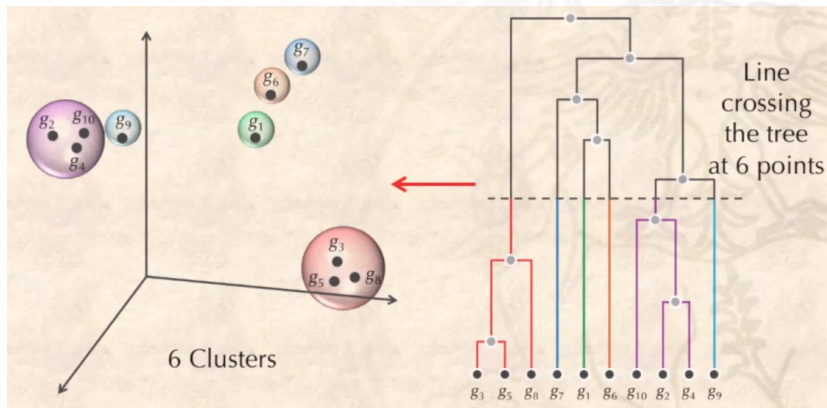
- To capture stratification, the **hierarchical clustering** algorithm organizes n data points into a tree.



Hierarchical Clustering

From Data to a Tree

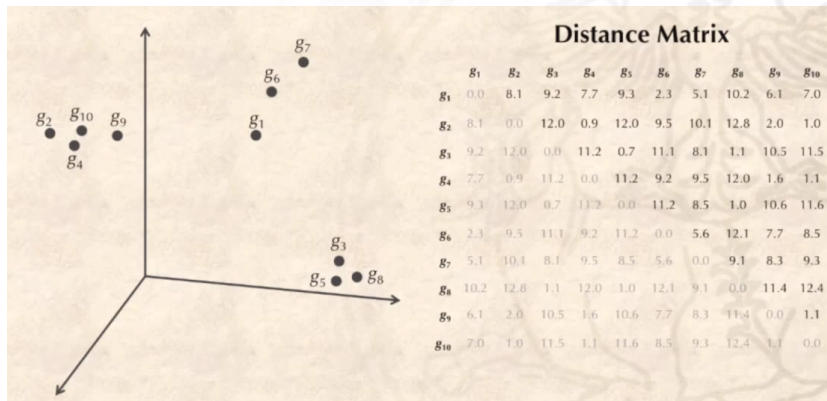
- To capture stratification, the **hierarchical clustering** algorithm organizes n data points into a tree.



Hierarchical Clustering

Constructing the Tree

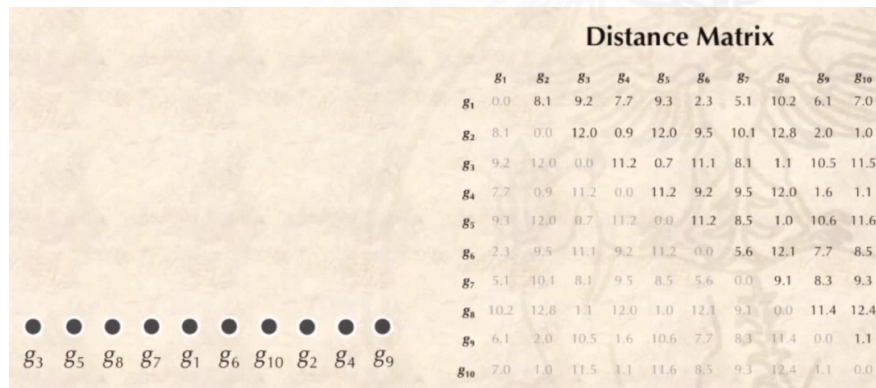
- Hierarchical clustering starts from a transformation of $n \times m$ expression matrix into $n \times n$ **similarity matrix** or **Distance matrix**;
- it can be obtained by simply computing Euclidean/Manhattan distance between genes.



Hierarchical Clustering

Constructing the Tree

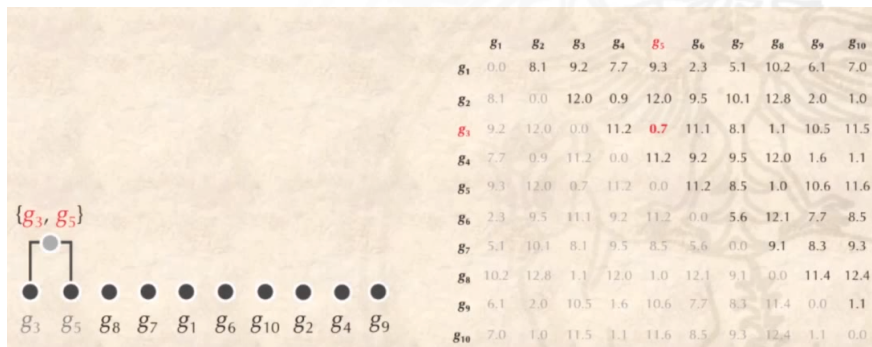
- Create a node (i.e. a single element cluster) for every gene.



Hierarchical Clustering

Constructing the Tree

- Identify the two *closest* clusters and merge them.

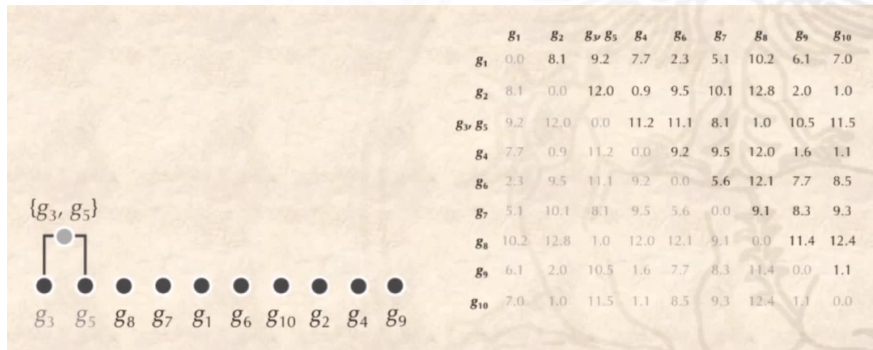


Hierarchical Clustering

Constructing the Tree

- Recompute the distance between two clusters as the minimal distance between the elements in the clusters:

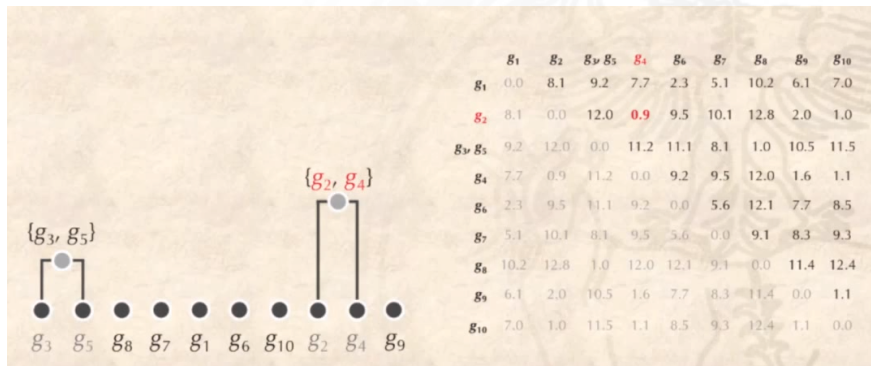
$$D(C_1, C_2) = \min_{\forall i \in C_1, j \in C_2} D_{i,j}$$



Hierarchical Clustering

Constructing the Tree

- Identify the two *closest* clusters and merge them.

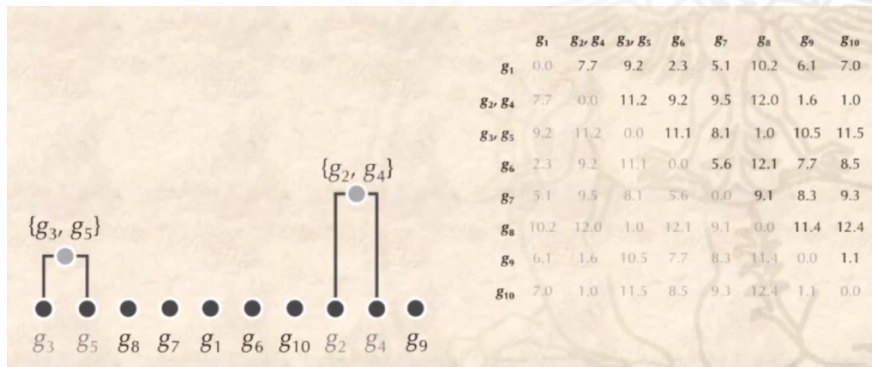


Hierarchical Clustering

Constructing the Tree

- Recompute the distance between two clusters as the minimal distance between the elements in the clusters:

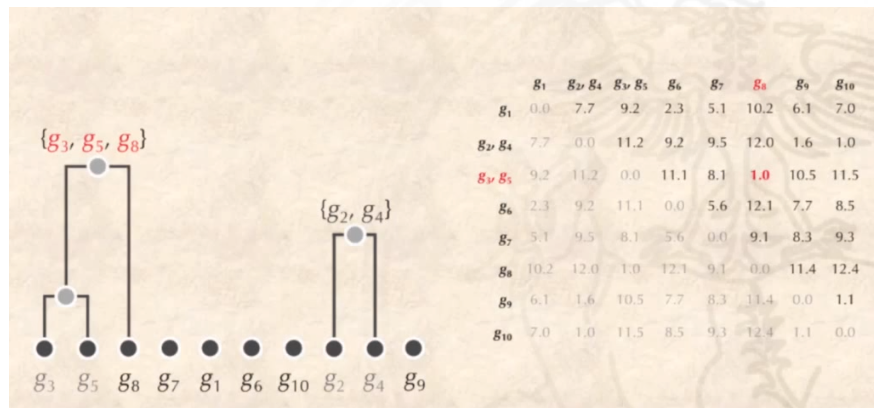
$$D(C_1, C_2) = \min_{\forall i \in C_1, j \in C_2} D_{i,j}$$



Hierarchical Clustering

Constructing the Tree

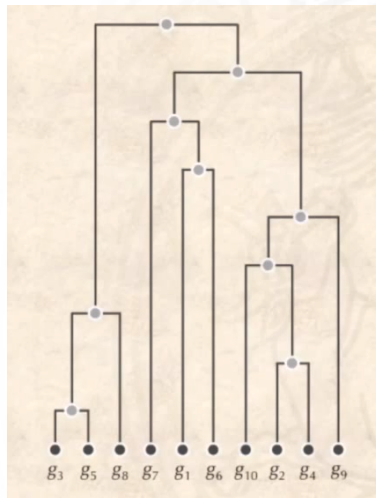
- Identify the two *closest* clusters and merge them.



Hierarchical Clustering

Constructing the Tree

- Iterate until all elements form a single cluster (i.e. *root*).



Hierarchical Clustering

Constructing the Tree

Hierarchical Clustering (D, n)

$Clusters \leftarrow n$ single-element clusters labeled 1 to n

$T \leftarrow$ a graph with the n isolated nodes labeled 1 to n

while there is more than one cluster

find the two closest clusters C_i and C_j

merge C_i and C_j into a new cluster C_{new} with $|C_i| + |C_j|$ elements

add a new node labeled by cluster C_{new} to T

connect node C_{new} to C_i and C_j by directed edges

remove the rows and columns of D corresponding to C_i and C_j

remove C_i and C_j from $Clusters$

add a row and column to D for the cluster C_{new} by computing

$D(C_{new}, C)$ for each cluster C in $Clusters$

add C_{new} to $Clusters$

assign root in T as a node with no incoming edges

return T

Hierarchical Clustering

Different distance functions

Minimum distance between elements of two clusters:

$$D_{\min}(C_1, C_2) = \min_{\text{all points } i \text{ and } j \text{ in clusters } C_1 \text{ and } C_2, \text{ respectively}} D_{i,j}$$

Average distance between elements of two clusters:

$$D_{\text{avg}}(C_1, C_2) = (\sum_{\text{all points } i \text{ and } j \text{ in clusters } C_1 \text{ and } C_2, \text{ respectively}} D_{i,j}) / (|C_1| * |C_2|)$$

Difference between K-means and Hierarchical clustering

- Hierarchical clustering can not handle big data as well as K-means clustering.
- In K-means clustering, since we start with random choice of clusters, the results produced by running the algorithm multiple times might differ. While results are reproducible in Hierarchical clustering;
- K-means clustering requires prior knowledge of K.

Exercises

- Try to apply K-means algorithm on the following inputs:
 - ▶ $k = 3$;
 - ▶ $P_1(1, 1), P_2(2, 1), P_3(3, 4), P_4(1, 6), P_5(2, 3), P_6(3, 3), P_7(4, 5), P_8(2, 5)$.
- Try to apply hierarchical clustering algorithm on the same input points.

$$MD = \begin{matrix} & g_1 & g_2 & g_3 & g_4 & g_5 & g_6 & g_7 & g_8 \\ g_1 & 0.000 & 1.000 & 3.605 & 5.000 & 2.236 & 2.828 & 5.000 & 4.123 \\ g_2 & 1.000 & 0.000 & 3.162 & 5.099 & 2.000 & 2.236 & 4.472 & 4.000 \\ g_3 & 3.605 & 3.162 & 0.000 & 2.828 & 1.414 & 1.000 & 1.414 & 1.414 \\ g_4 & 5.000 & 5.099 & 2.828 & 0.000 & 3.162 & 3.605 & 3.162 & 1.414 \\ g_5 & 2.236 & 2.000 & 1.414 & 3.162 & 0.000 & 1.000 & 2.828 & 2.000 \\ g_6 & 2.828 & 2.236 & 1.000 & 3.605 & 1.000 & 0.000 & 2.236 & 2.236 \\ g_7 & 5.000 & 4.472 & 1.414 & 3.162 & 2.828 & 2.000 & 0.000 & 2.000 \\ g_8 & 4.123 & 4.000 & 1.414 & 1.414 & 2.000 & 2.236 & 2.000 & 0.000 \end{matrix}$$