# QSAR Toolbox functionalities.
# Clustering

Laboratory of  Mathematical Chemistry
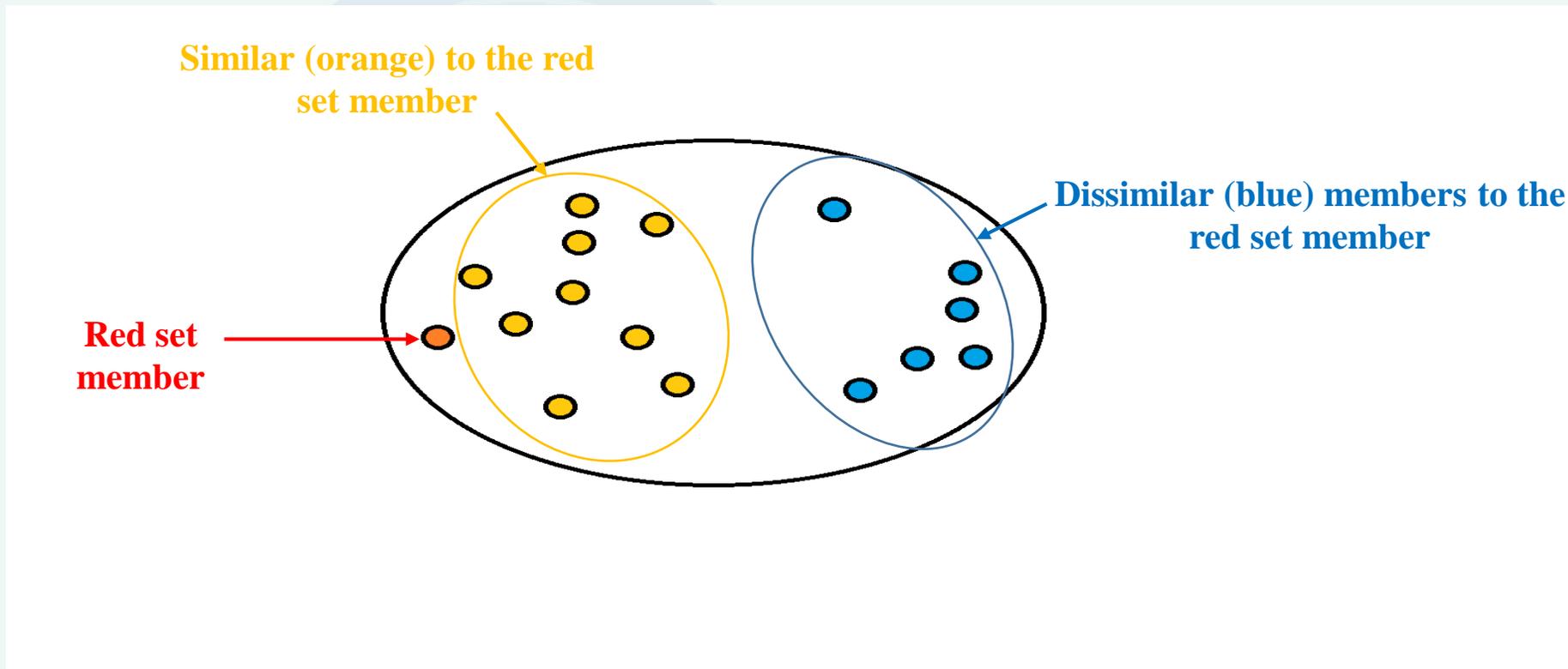University "Prof. As. Zlatarov" Bourgas, Bulgaria

June 2019

# Clustering - overview

o    Cluster analysis (clustering) is the task of grouping a set of objects in such a way that objects in the same group (called cluster) are more similar (in some sense or another) to each other than to those in other groups (clusters) [1].

o    Cluster analysis is based on similarities between each two chemicals of an initial set.  Similarity values are compared to a selected threshold in a repeated procedure.

o    Resulting clusters include only substances which are similar to each other above the required threshold.

[1] - http://en.wikipedia.org/wiki/Cluster_analysis

# Clustering – how the clusters are formed?

1) The system prepares a matrix which contains the similarities between each two items of the initial set. Another matrix is also prepared that contains only flags if each two items are considered "similar" or "dissimilar". The latter values are calculated by comparing the similarity values from the first matrix to a given threshold.
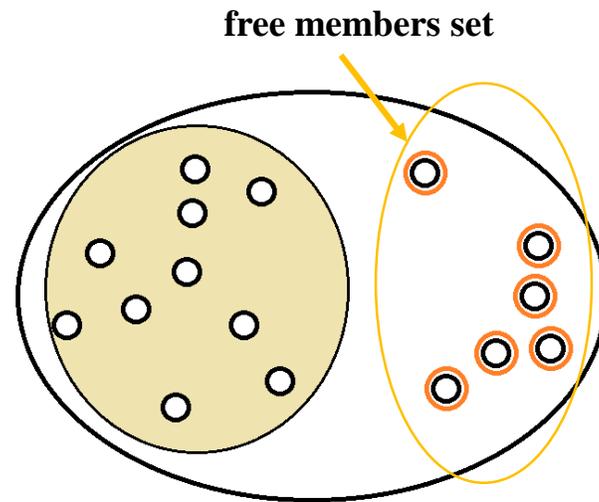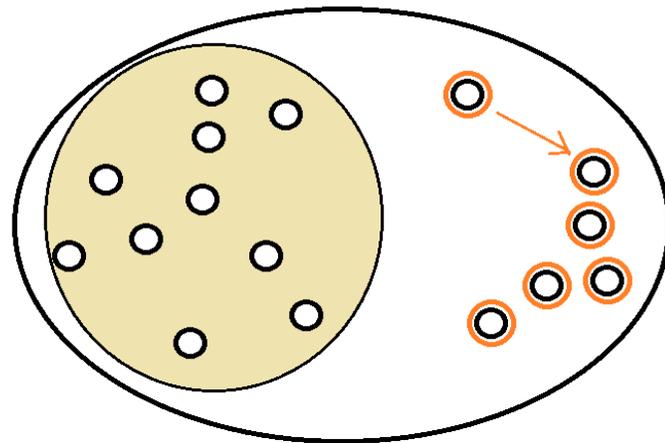
# Clustering – how the clusters are formed?

2) After creating a cluster, if there are free members set (not yet included in the existing clusters) still remaining, the process is repeated again.

Initially each free member* breeds a new cluster. Then, every cluster grows as much as possible, including only members that are similar to the ones it already included.

Depending on the option "***Allow overlapping***" the cluster may or may not include members that are present in already existing clusters.
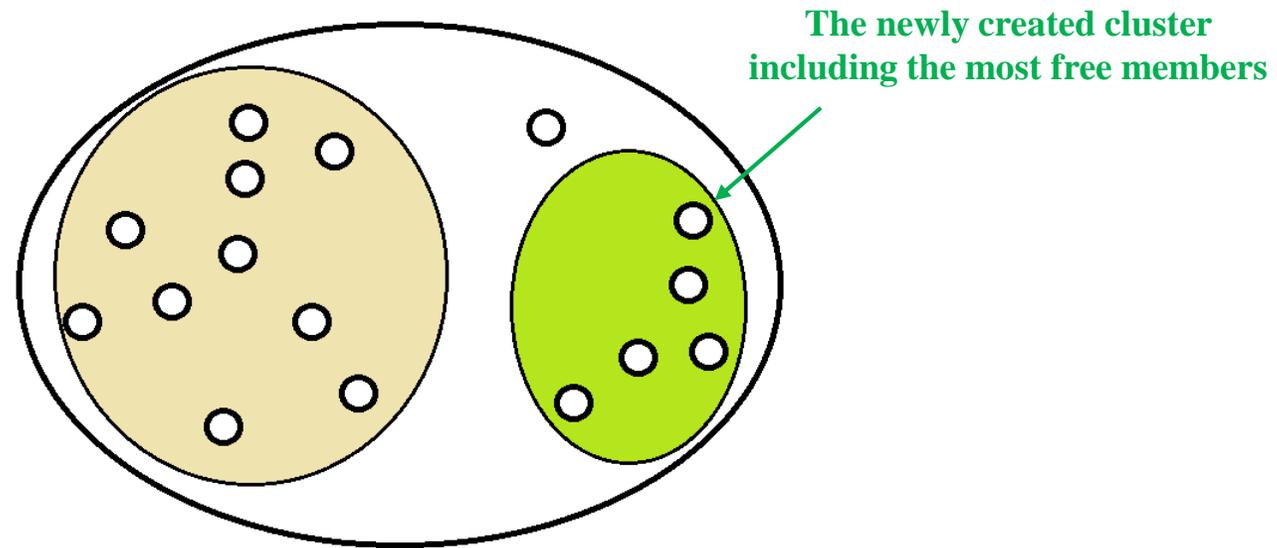


free members set

*Always checked is the nearest free available member. The distance between a cluster and a free member is measured through the distances between this free member and the members already included in the new cluster. Different options are available here – ***taking the minimal***, ***taking the maximal***, ***taking the average***, ***taking the median***.

# Clustering – how the clusters are formed?

2) After creating a cluster, if there are free members set (not yet included in the existing clusters) still remaining, the process is repeated again.

Initially each free member* breeds a new cluster. Then, every cluster grows as much as possible, including only members that are similar to the ones it already included.

Depending on the option "***Allow overlapping***" the cluster may or may not include members that are present in already existing clusters.

**Every free member breeds a new cluster; the cluster grows to the nearest member**



*Always checked is the nearest free available member. The distance between a cluster and a free member is measured through the distances between this free member and the members already included in the new cluster. Different options are available here – ***taking the minimal***, ***taking the maximal***, ***taking the average***, ***taking the median***.

# Clustering – how the clusters are formed?

3) At the end, selected is the cluster which had included most free members

**The newly created cluster including the most free members**

# Clustering – implementation in QSAR Toolbox

The clustering functionality in Toolbox is available within the Category definition module.

# Clustering – implementation in QSAR Toolbox

Allows distributing defined category into clusters. Cluster is presented as a sub-category of general category that includes chemicals with unique combination of profiling results or similar structural characteristics*.



***Note**: When *Structure similarity* profiler is selected to cluster a category, the sub-categories depends on the defined threshold and target chemical availability.

    1) If there is a defined target chemical, then the category is splitted into clusters based on the structural similarity of the chemicals with respect to the target

    2) If there is not a defined target chemical, then the category is splitted into clusters based on the structural similarity between each to each chemical

# Clustering – implementation in QSAR Toolbox

Allows distributing defined category into clusters. Cluster is presented as a sub-category of general category that includes chemicals with unique combination of profiling results or similar structural characteristics.



**All chemicals in one cluster have the same profiling results according to the selected profiler**