

Homework 2

Pattern Mining and Social Network Analysis

BOUYSSOU Gatien , de POURTALES Caroline, LAMBA Ankit

10 novembre, 2020

Contents

Cluster validation techniques	3
Clustering Tendency	3
Why assessing clustering tendency ?	3
Methods	3
Statistical methods for clustering tendency	3
Visually	3
Statistics on a model	3
Internal measures	3
Compactness	4
Separation	4
Connectivity	4
Average silhouette	4
Silhouette coefficient of an observation	4
Average silhouette	4
Dunn index	4
External mesures	4
Determining the Optimal number of Clusters	5
Elbow method	5
Average silhouette method	5
Gap statistic method	5
Principal Components Analysis	6
Different kinds of PCA	6
Standard PCA	6
Deciding how many PCs to use	6
Proportion of variance explained (PVE)	6
Cumulative Proportion of variance explained	6
Incremental PCA	6
Sparse PCA	6
Kernel PCA	7
Example on tissue samples	7
On R	7
On Python with scikit-learn	8
Clustering	9
Partitioning Clustering	9
K-means	9
Within-cluster variation (squared Euclidean distance)	9
K-means algorithm	9

Choice of k	9
Example on tissue samples	9
On R	9
On Python with scikit-learn	10
K-medoids algorithm	11
Principle	11
PAM algorithm (Partitioning Around Medoids)	11
CLARA - Clustering Large Applications	12
Hierarchical clustering	13
Dissimilarity function	13
Euclidean distance	13
Correlation-based distance	13
Linkage	13
Maximum or complete linkage	13
Minimum or single linkage	13
Mean or average linkage	13
Centroid linkage	13
Ward's minimum variance method	13
P-value	14
Examples	14
Example on tissues samples	14
On R	14
On Python with scikit-learn	15
Example on Corona dataset : the problem of over-represented category	16
On R	17
On Python with scikit-learn	17
Fuzzy Clustering	18
Overall	18
Examples	18
Choosing the best algorithm ?	20
Four Basic Algorithms And How To Choose One	20
Measures for comparing clustering algorithms	20

Cluster validation techniques

Clustering Tendency

Why assessing clustering tendency ?

Clustering can create clusters even if there are no meaningful cluster.

Clustering tendency assessment methods are used to evaluate the validity of clustering analysis. It means evaluating if there are meaningful clusters in the data.

Methods

Statistical methods for clustering tendency A method called Hopkins statistic is used to assess the clustering tendency in a data set. It measures the probability for the data set to be generated by an uniform data distribution, it means that this statistic looks at the spatial randomness of the data.

For D , a data set :

- Sample n points from D (p_1, \dots, p_n)
 - For all $p_i \in D$, compute the distance to the nearest neighbour x_i
- Generate ($random_D$) from a random uniform distribution with n points (q_1, \dots, q_n) with the same variation as D .
 - For all $q_i \in random_D$, compute the distance to the nearest neighbour y_i
- Calculate the Hopkins statistic (H). It is the mean of the nearest neighbour distances in $random_D$ divided by the sum of the mean nearest neighbour distances in D and $random_D$.

$$H = \frac{\sum_i^n y_i}{\sum_i^n x_i + \sum_i^n y_i}$$

If H is about 0.5, it means that the sum along D and $random_D$ are very close so that the data D is uniformly distributed (no meaningful clusters).

It is called the null hypothesis.

Otherwise, we have the Alternative hypothesis: the data set D is not uniformly distributed (contains meaningful clusters).

Visually There is an algorithm for visual assessment of cluster tendency (VAT) : it is the approach by Bezdek and Hathaway, 2002.

- Compute the dissimilarity matrix using the Euclidean distance measure
- Reorder the DM so that similar objects are close to one another. It is now called the ODM.
- Display the ODM as an ordered dissimilarity image (ODI) (with some libraries).

Computing the visual form gives us colored squares along the diagonal. The VAT detects the clustering tendency in a visual form by counting them.

Statistics on a model

Internal measures

Internal cluster validation uses the internal information of the clustering process to evaluate the clustering structure without any external information.

Internal validation measures reflect often the compactness, the connectedness and the separation of the cluster partitions.

Compactness Compactness measures how close the objects are, within the same cluster. We evaluate it with the notion of distance such as the cluster-wise within average/median distances between observations.

Separation Separation measures how well clusters are separated from one another. We evaluate it by looking at the distances between clusters' centers or with the pairwise minimum distances between objects in different clusters.

Connectivity Connectivity corresponds to what extent items are placed in the same cluster as their nearest neighbours in the data space.

Average silhouette

Silhouette coefficient of an observation It measures how well an observation is clustered and it estimates the average distance between clusters.

It is possible to compute this coefficient thanks to the following formula :

$$Silhouette_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

Where :

- a_i is the average distance between the observation i within its cluster
- b_i is the average distance between the observation i and all the observations belonging to another cluster

How to interpret the value of S_i ?

A large S_i , almost 1, means the observation is very well clustered.

A small S_i , close to 0, means that the observation lies between two clusters.

A negative S_i means the observation is probably placed in the wrong cluster.

Average silhouette The average silhouette is the mean of all silhouette coefficients. It is used to evaluate the result made by a clustering algorithm.

Dunn index The Dunn index is another internal clustering validation measure.

It is the minimal distance between two clusters (the smallest separation) over the maximal distance between the objects of one clusters (the biggest diameter of a cluster).

$$Dunn\ index = \frac{min.separation}{max.diameter}$$

Compact and well-separated clusters in a data set means a small diameter of the clusters and a large distance between the clusters. Thus, Dunn index should be maximized.

External mesures

External cluster validation compares the results of a cluster analysis with an externally known result. It measures how well a cluster match extern class labels. We know k , the number of clusters, in advance. So we use this validation to choose the correct clustering method.

We compare the identified clusters to an external reference.

Determining the Optimal number of Clusters

Determining the optimal number of clusters in a data set is fundamental because, for example, in partitioning clustering, such as k-means clustering, it requires the user to specify the number of clusters k .

Elbow method

The Elbow method looks at the total within-cluster sum of squares over the number of clusters.

- Compute the clustering algorithm chosen with k going from 1 to 10 for example.
- Calculate the within-cluster sum of squares for each k
- Plot the curve of WSS according to k
- The location of a bend (knee) in the plot is generally considered as an indicator of the appropriate number of clusters.

An alternative is the Silhouette method.

Average silhouette method

Average silhouette method computes the average silhouette of observations for different values of k . Then, k is chosen if it maximizes the average silhouette.

- Compute the clustering algorithm chosen with k going from 1 to 10 for example.
- Calculate the average silhouette coefficients of observations for each k
- Plot the curve of AVG.S according to k
- The location of the maximum is considered as the appropriate number of clusters.

Gap statistic method

The gap statistic compares the total within intra-cluster variation for different values of k with their expected values under null reference. We choose k when it maximizes the gap statistic so that the clustering structure is far away from the random uniform distribution of points.

- Compute the clustering algorithm chosen with k going from 1 to 10 for example.
- Calculate the within-cluster sum of squares for each k W_k
- Generate B reference data sets with a random uniform distribution.
 - Compute the clustering algorithm chosen on these reference data sets with k going from 1 to k_m *ax*.
 - Calculate the within-cluster sum of squares for each k W_{kb}
- Compute the gap :

$$Gap(k) = \frac{1}{B} \sum_{b=1}^B \log(W_{kb}) - \log(W_k)$$

- Compute the standard deviation of the statistics s_k .
- Choose the smallest k such that the gap statistic is within one standard deviation of the gap at $k+1$.

$$Gap(k+1) > Gap(k) - s_{k+1}$$

Principal Components Analysis

The goal of PCA is to identify which features, in the dataset, explain the most variability.

Principal Component Analysis, or PCA, is a dimensionality-reduction method that is often used to reduce the dimensionality of large data sets, by transforming a large set of variables into a small one that still contains most of the information from the large set.

Different kinds of PCA

Standard PCA

To process the feature extraction, we need to associate eigenvectors and eigenvalues to each feature of our dataset. For that, we will follow those steps one by one :

- First, put your class label column aside and keep all the independent features (the rest)
- For each remaining columns subtract the mean of that column from each entry
- Standardize the features if needed. This will form a new matrix S
- Transpose this matrix and multiply it by itself ($S^T S$). This is a new matrix.
- Compute the eigenvectors and the eigenvalues of this new matrix and deduce its eigendecomposition PDP^{-1} where P contains all the eigenvectors and D all the eigenvalues on the diagonal. Each eigenvalue is associated to a eigenvector (the first eigenvalue is associated to first column).
- Sort the eigenvalues. Change P based on the new order given to the eigenvalues. This new matrix will be called P^*
- Then multiply S by P^* . SP^* is a centered and standardized version of independent features.

Now that we have all the eigenvectors and eigenvalues corresponding to our features, we can determine how much features we want to keep.

Deciding how many PCs to use

For this part, we are going through two different methods the Proportion of variance explained and the Cumulative Proportion of variance explained.

Proportion of variance explained (PVE)

This value for each column can be obtained by taking the corresponding eigenvalue on the diagonal matrix and dividing it by the sum of all the eigenvalues. The graph representing the proportion of variance explained should be a decreasing curve since the eigenvalues have been sorted from the greatest to the lowest.

Cumulative Proportion of variance explained

The Cumulative Proportion of Variance Explained (CPVE) is the sum of the first components divided by all the others. For example, the CPVE for the third component is

$$CPVE = \frac{e_1 + e_2 + e_3}{e_1 + e_2 + e_3 + \dots e_k}$$

Where e_i is the eigenvalue for the i th value and k is the number of columns

Incremental PCA

The IncrementalPCA object uses a different form of processing and allows for partial computations which almost exactly match the results of PCA while processing the data in a minibatch fashion.

Sparse PCA

SparsePCA is a variant of PCA, with the goal of extracting the set of sparse components that best reconstruct the data.

Kernel PCA

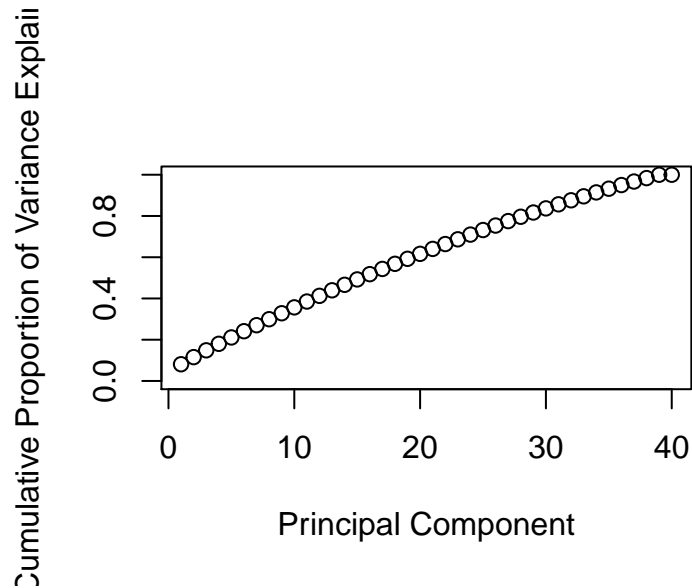
KernelPCA is an extension of PCA which achieves non-linear dimensionality reduction through the use of kernels.

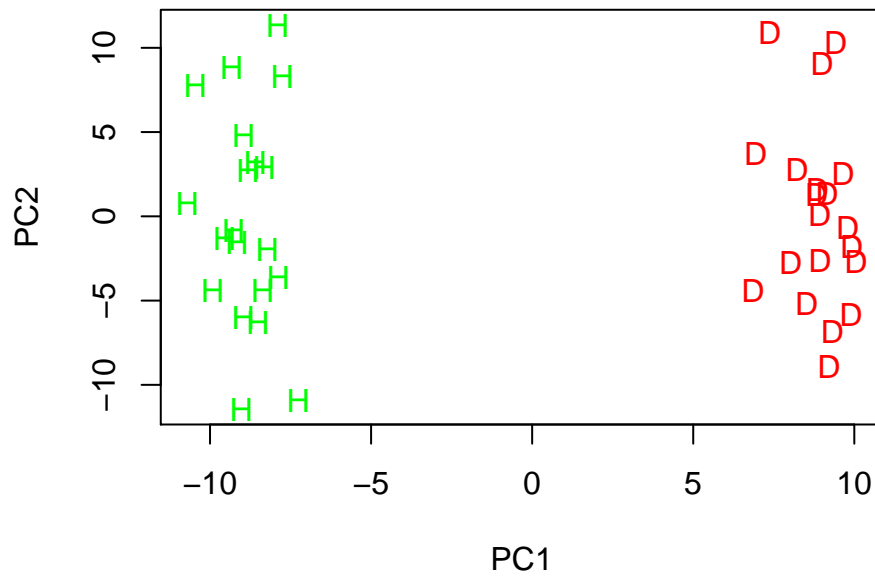
Example on tissue samples

The following dataset consists of 40 tissue samples with measurements of 1,000 genes. The first 20 tissues come from healthy patients (H) and the remaining 20 come from a diseased patient group (D).

```
## 40
## False
## <class 'pandas.core.frame.DataFrame'>
## RangeIndex: 40 entries, 0 to 39
## Columns: 1000 entries, 0 to 999
## dtypes: float64(1000)
## memory usage: 312.6 KB
```

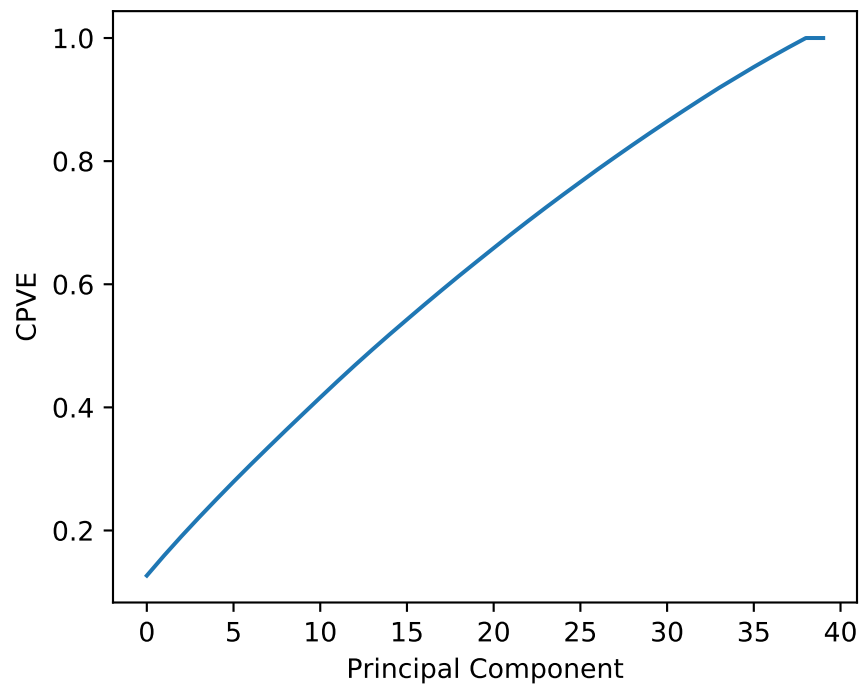
On R





On Python with scikit-learn

```
## PCA(n_components=40)
```



According to the Cumulative Proportion of Variance Explained (CPVE), we can see that the first 35 eigenvalues correspond to approximately 98% of the variance. It means that we need only 35 patients on 40.

Clustering

Partitioning Clustering

K-means

The objective of clustering is to distinct groups from the datatest. With k-means, we want to distinct k groups. The algorithm will assign each observation to exactly one of the cluster. It optimizes the groups by minimizing the within-cluster variation such that the sum of the with-cluster variations across all the clusters is the smallest possible.

Within-cluster variation (squared Euclidean distance) If μ_k is the center of the cluster k. The total with-cluster variation is TW :

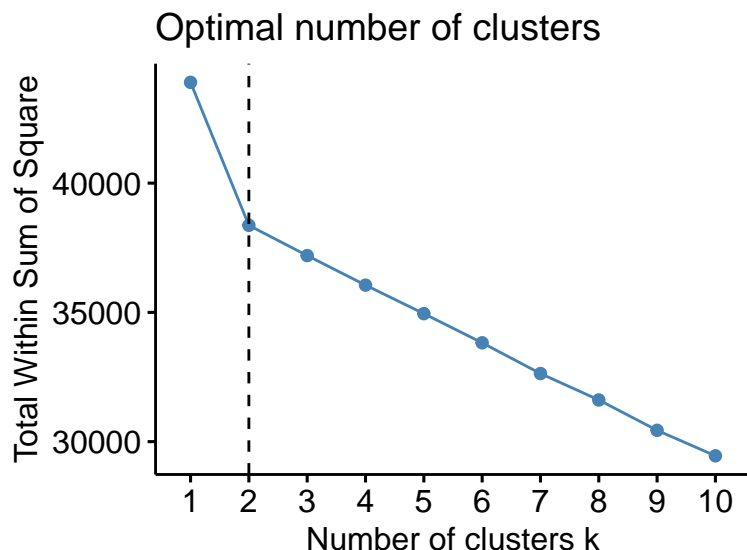
$$TW = \sum_{j=1}^k W_j = \sum_{j=1}^k \sum_{x_i \in C_j} (x_i - \mu_k)^2$$

K-means algorithm The first step when using k-means clustering is to indicate the number of clusters (k) that will be generated in the final solution. The algorithm starts by randomly selecting k objects from the data set to serve as the initial centers for the clusters. The selected objects are also known as cluster means or centroids.

Choice of k We compute k-means clustering using different k, then we choose the number of cluster according to the location of a bend on the graph representing the Within-cluster variation according to k.

Example on tissue samples

On R According to this graph, we should choose k=2 (it makes sense since we have Healthy and non healthy patients).



Then, applying kmeans with 2 clusters we observe that the 20 first individuals (healthy) are not in the same cluster than the 20 others (non healthy).

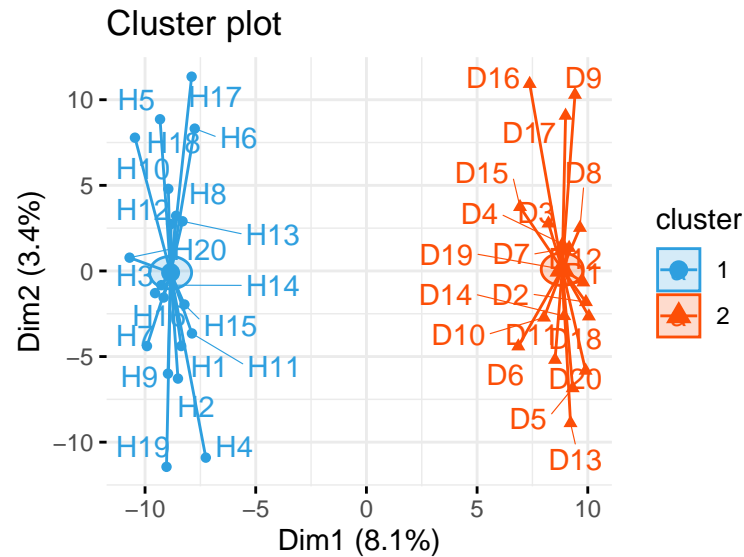
##	H1	H2	H3	H4	H5	H6	H7	H8	H9	H10	H11	H12	H13	H14	H15	H16	H17	H18	H19	H20
##	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
##	D1	D2	D3	D4	D5	D6	D7	D8	D9	D10	D11	D12	D13	D14	D15	D16	D17	D18	D19	D20
##	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2

The following coefficient is the Dunn index.

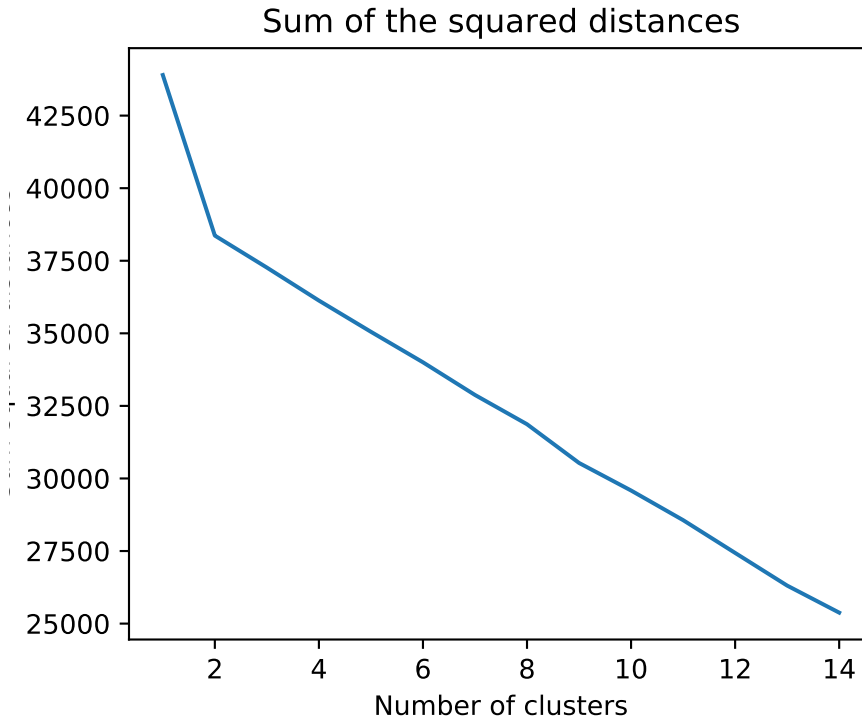
It is about 1. It means that the separation distance between the two clusters is almost the same as the diameter of the largest cluster.

```
## [1] 0.9744956
```

Since we have a multi-dimensional dataset, we apply dimensionality reduction with the use of PCA to plot the clusters. On the x axis, it is the first PCA, on the y axis, it is the second PCA.



On Python with scikit-learn With python, we are going to repeat the experiment to validate the number of cluster and the different results found above. At first, we are going to loop through the number of clusters k, and select one, based on the sum of squared distances.



Using the values and the graph, we can determine using the elbow method that 2 is the optimal number of clusters since it is the spot where the curve is bending.

This is the silhouette score :

```
## 0.09607065301278583
```

The silhouette coefficient estimates the average distance between clusters. If the coefficient tends toward 1, it means that the values are well clustered. In that case, the silhouette score is equal to 0. It means that the clusters are overlapping. However, the silhouette score is not negative. It would have meant that the sample has been assigned to the wrong cluster.

This is the dunn index :

```
## 2.633190807264714
```

The dunn index strengthens the idea that the values are well classified. If they were not well classified the dunn index would have been close to 0.

K-medoids algorithm

Principle The k-medoids algorithm is a clustering approach related to k-means clustering. In k-medoids clustering, each cluster is represented by one of the data point in the cluster.

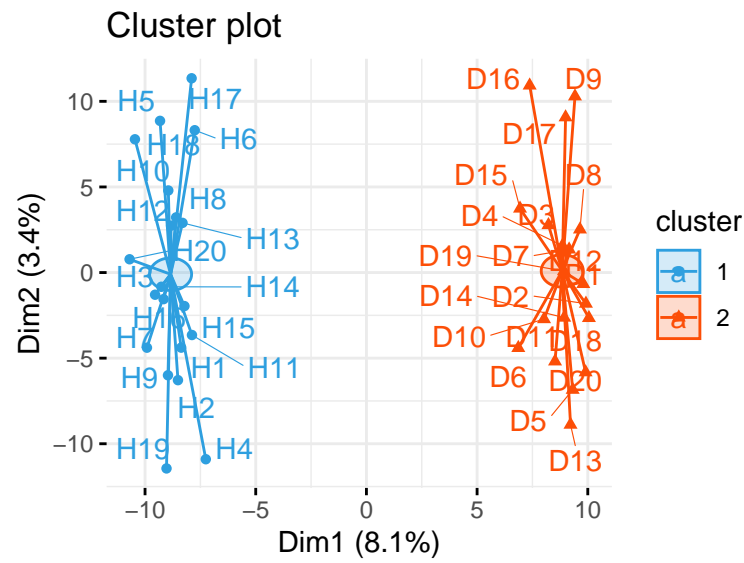
The most common k-medoids clustering methods is the PAM.

PAM algorithm (Partitioning Around Medoids)

```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

We have the same Dunn index as with K-means on the same dataset. Both methods perform equally.

[1] 0.9744956



CLARA - Clustering Large Applications CLARA (Clustering Large Applications, (Kaufman and Rousseeuw 1990)) is an extension to k-medoids (PAM) methods to deal with data containing a large number of objects (more than several thousand observations) in order to reduce computing time and RAM storage problem. This is achieved using the sampling approach.

Hierarchical clustering

Dissimilarity function

Euclidean distance

##		H1	H2	H3	H4	H5	H6	H7	H8	H9	H10
##	H1	0.0	45.2	44.5	45.6	45.9	45.2	45.0	45.8	45.0	46.0
##	H2	45.2	0.0	45.7	44.6	44.4	45.5	44.8	44.6	44.4	44.5
##	H3	44.5	45.7	0.0	43.9	46.0	44.9	45.7	45.2	44.1	45.3
##	H4	45.6	44.6	43.9	0.0	47.4	45.1	45.6	45.7	44.0	44.4
##	H5	45.9	44.4	46.0	47.4	0.0	45.4	44.6	45.6	45.7	44.3
##	H6	45.2	45.5	44.9	45.1	45.4	0.0	44.9	43.9	44.7	43.3
##	H7	45.0	44.8	45.7	45.6	44.6	44.9	0.0	45.3	43.5	44.4
##	H8	45.8	44.6	45.2	45.7	45.6	43.9	45.3	0.0	44.0	44.1
##	H9	45.0	44.4	44.1	44.0	45.7	44.7	43.5	44.0	0.0	44.2
##	H10	46.0	44.5	45.3	44.4	44.3	43.3	44.4	44.1	44.2	0.0

Correlation-based distance Correlation-based distance considers two observations to be similar if their features are highly correlated, even though the observed values may be far apart in terms of Euclidean distance.

##		H1	H2	H3	H4	H5	H6	H7	H8	H9	H10
##	H1	0.0	1	1.0	1.0	1.0	1.0	1.0	1.1	1.0	1.1
##	H2	1.0	0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
##	H3	1.0	1	0.0	0.9	1.0	1.0	1.0	1.0	1.0	1.0
##	H4	1.0	1	0.9	0.0	1.1	1.0	1.0	1.0	1.0	1.0
##	H5	1.0	1	1.0	1.1	0.0	1.0	1.0	1.0	1.0	1.0
##	H6	1.0	1	1.0	1.0	1.0	0.0	1.0	1.0	1.0	0.9
##	H7	1.0	1	1.0	1.0	1.0	1.0	0.0	1.0	0.9	1.0
##	H8	1.1	1	1.0	1.0	1.0	1.0	1.0	0.0	1.0	1.0
##	H9	1.0	1	1.0	1.0	1.0	1.0	0.9	1.0	0.0	1.0
##	H10	1.1	1	1.0	1.0	1.0	0.9	1.0	1.0	1.0	0.0

Linkage

The linkage function takes the distances and groups pairs of objects into clusters based on their similarity. These clusters are then linked to each other to create bigger clusters and the linkage continues until all the data are linked together in a hierarchical tree.

Maximum or complete linkage The distance between two clusters is defined as the maximum value of all pairwise distances between the elements in cluster 1 and the elements in cluster 2. It tends to produce more compact clusters.

Minimum or single linkage The distance between two clusters is defined as the minimum value of all pairwise distances between the elements in cluster 1 and the elements in cluster 2. It tends to produce long, “loose” clusters.

Mean or average linkage The distance between two clusters is defined as the average distance between the elements in cluster 1 and the elements in cluster 2.

Centroid linkage The distance between two clusters is defined as the distance between the centroid for cluster 1 (a mean vector of length p variables) and the centroid for cluster 2.

Ward’s minimum variance method It minimizes the total within-cluster variance. At each step the pair of clusters with minimum between-cluster distance are merged.

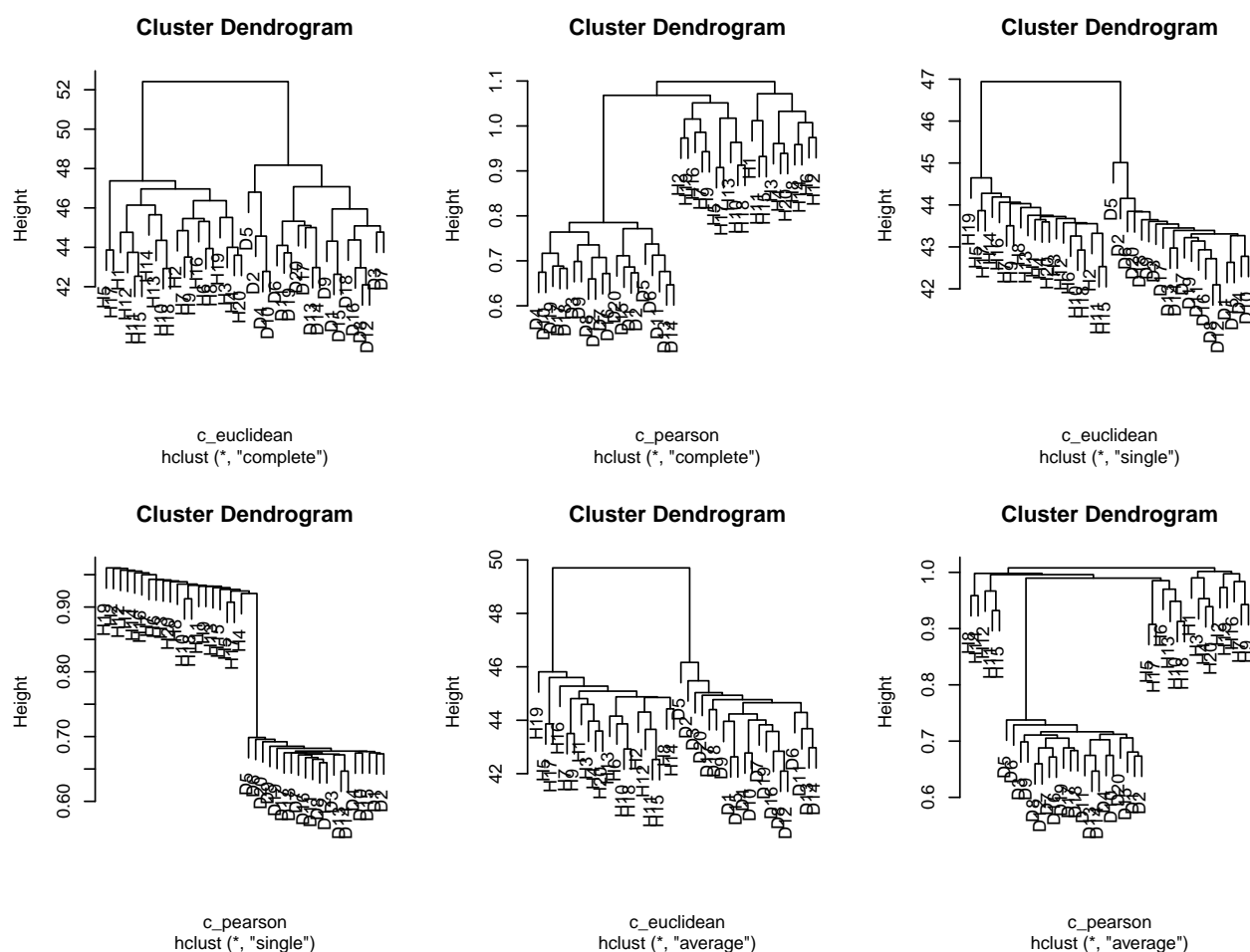
P-value

- Generate bootstraps samples
- Compute hierarchical clustering on each bootstrap copy
- Compute for each cluster :
 - The bootstrap probability (BP) value : frequency that the cluster is identified in bootstrap copies.
 - The approximately unbiased (AU) probability values (p-values) by multiscale bootstrap resampling. Clusters with p-value above 95% are considered to be strongly supported by data.

Examples

Example on tissues samples Using the tissue and genes data set, we apply hierarchical clustering.

On R



We can see that the use of euclidean distance in the three methods (complete, single, average) gives good results (no missclassification) but the use of correlation-distance gives very bad results.

Furthermore all methods, except Average with correlation-distance, divide the graph in two groups (healthy and non-healthy) which is very good.

We have the same Dunn index as with K-means and PAM on the same dataset (tissues samples). So k-means, PAM and hierarchical clustering (with euclidean distance and complete linkage) methods perform equally.

```
## [1] 0.9744956
```

The cluster average silhouette widths are :

```
##          1          2
## 0.09268309 0.09945822
```

It means that observations are well clustered but the distinction between clusters is not that easy.

On Python with scikit-learn Thanks to hierarchical clustering we will be able to remove some useless features and to simplify the dataset. This will allow us to easily display the data thanks to a graph.

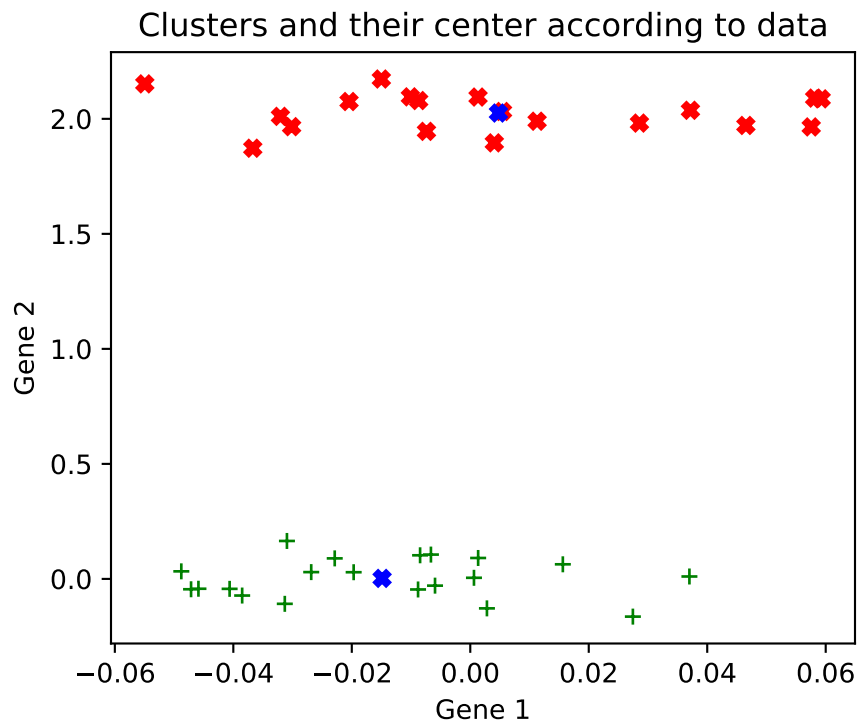
Feature Agglomeration begin with each element as its own cluster and recursively merges them to make larger clusters. In argument, it can take a linkage (complete, ward, single, average) and an affinity as distance.

```
## FeatureAgglomeration()
## (40, 2)
```

This hierarchical clustering allows us to keep only 2 features over 1000. It means that only 2 features were determinant to classify the 40 patients.

With the dataset simplified it will be easier to plot the performances of the Kmeans algorithm (with only 2 genes).

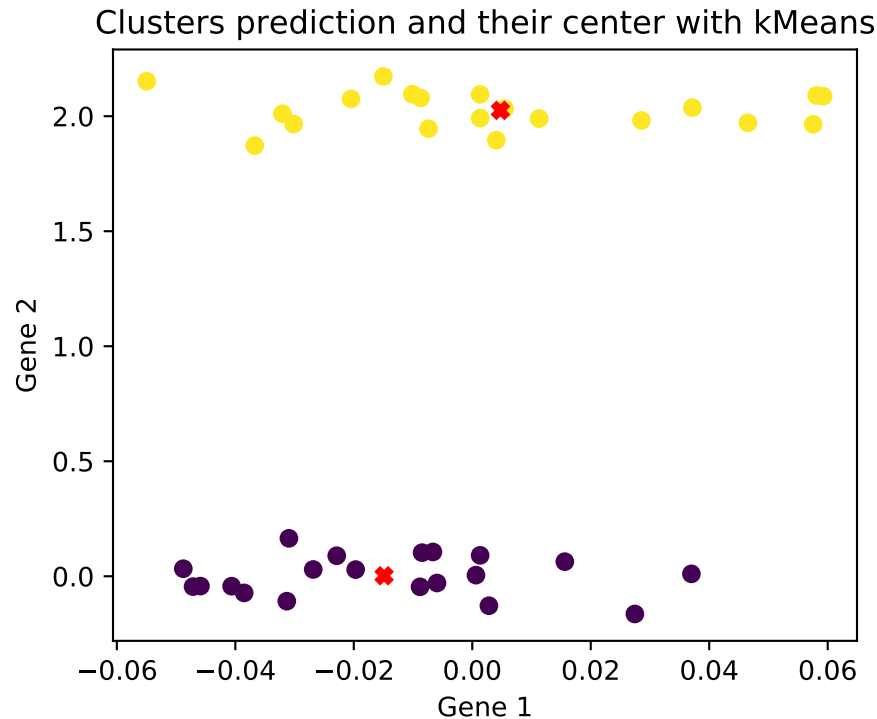
At first, we plot the points with the colors changing depending on the real class label.



The clustering obtained is :

```
## array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1,
##        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1], dtype=int32)
```

Then, we are plotting the points with the colors changing depending on the predicted class label.



As we can see the predicted class correspond to the real class labels. The clusters are well done.

Example on Corona dataset : the problem of over-represented category The corona dataset contains information of patients who have caught Corona. It stores their age, sex and nationality but also if there are dead or not.

The purpous is to use hiarchical clustering to see if the age, sex and nationality is enough to separate the dead from the living patients. H stands for Healthy now and D stands for Deceased.

```
##      deceased      sex  age country
## H1          0   male   78  France
## H2          0   male   94  France
## H3          0 female   87  France
## H4          0 female   66  France
## H5          0   male   32  France
```

We can observe here that the columns sex and country contain strings. Therefore, we need to convert those strings into numbers.

This is the percentage of death in the data set :

```
## 9.090909090909092
```

As we can see the dataset is unbalanced. There is only 2% of positive case for 98% negative case in the database.

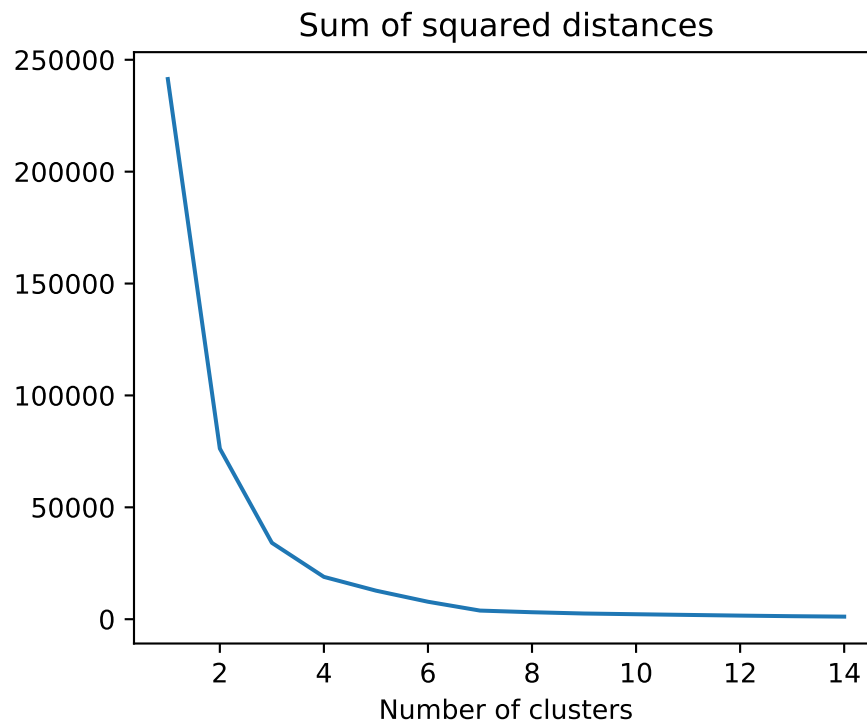
On R Since we have discrete features, we use the daisy distance.

```
## H446 H447 H448 H449 H450 D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D12 D13 D14 D15 D16 D17 D18 D19 D20 D21 D22 D23 D24 D25 D26 D27
## 1 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1
## D28 D29 D30 D31 D32 D33 D34 D35 D36 D37 D38 D39 D40 D41 D42 D43
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D44 D45
## 1 1
```

Looking at the clustering, 1 is the cluster of healthy patients, 2 is the cluster of dead patients. We can see that it is difficult to cluster the dead patients together because there are only 45 deceased cases against 450 healthy patients. This is a problem of over-represented category.

On Python with scikit-learn In python we are going to find out that the dataset is unbalanced. The aim is to show how an unbalanced dataset influence the results of the clustering algorithms.

Now, the dataset is ready to be processed by our clustering algorithms.



First of all, we can see on this graph that it is not that simple to determine the optimal number of clusters. We see a bend around 2 so it should be 2. Moreover, there is only 2 different labels : deceased = 1 or deceased = 0.

```
## 0.5797979797979798
```

We can see that because the dataset is not balanced the accuracy score is low.

The silhouette score is :

```
## 0.6023960067628296
```

The dunn index is :

```
## 0.5373248073274328
```

According to the silhouette score the average distance between clusters is small but better than the tissue dataset. Moreover, thanks to the dunn index we can conclude that the diameter of the cluster is 2 times bigger than the separation intra clusters.

Fuzzy Clustering

Overall

Fuzzy clustering is a clustering method where data points can belong in more than one group ("cluster"). Clustering divides data points into groups based in similarity between items and looks to find patterns or similarity between items in a set; Items in clusters should be as similar as possible to each other and as dissimilar as possible to items in other groups. Computationally, it's much easier to create fuzzy boundaries than it is to settle on one cluster for one point.

Examples

The function `fanny()` computes fuzzy clustering.

We use the tissue dataset.

For the dissimilarity function, the choice was euclidean distance.

2 is for the number of clusters expected.

stand indicates whether variables are standardized before calculating the dissimilarities.

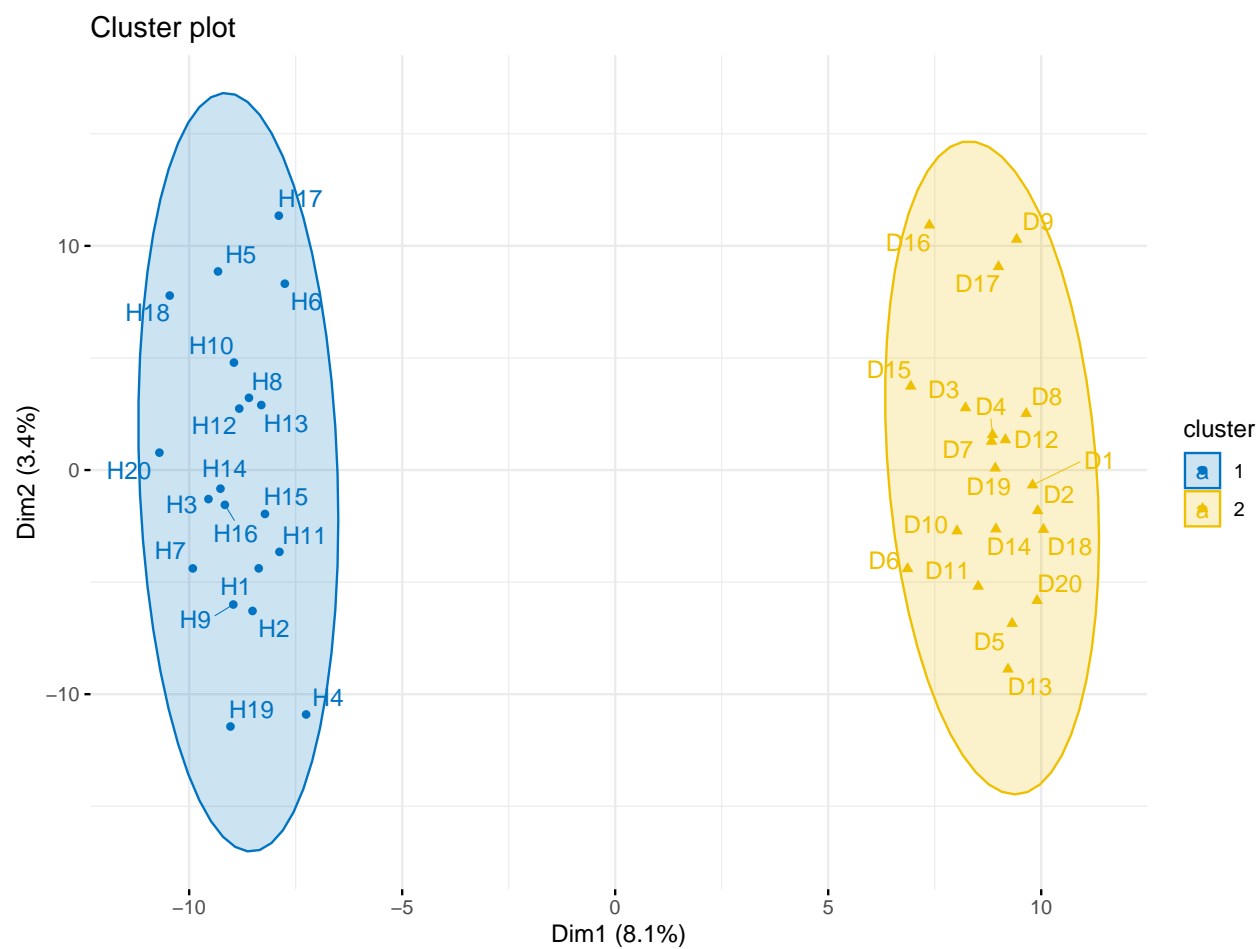
```
## H1 H2 H3 H4 H5 H6 H7 H8 H9 H10 H11 H12 H13 H14 H15 H16 H17 H18 H19 H20
## 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## D1 D2 D3 D4 D5 D6 D7 D8 D9 D10 D11 D12 D13 D14 D15 D16 D17 D18 D19 D20
## 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

We can see that every tissue sample was well clustered.

```
## dunn_coeff normalized
## 5.000000e-01 2.220446e-15
```

Looking at the Dunn coefficient, it is about 0.5 meaning that the diameter of one of the cluster is 2 times larger than the minimal separation of two cluster.

This makes sense looking at the following plot.



Choosing the best algorithm ?

Applying a clustering algorithm is much easier than selecting the best one. Each type offers pros and cons that must be considered if you're striving for a tidy cluster structure.

Clustering is an essential step in the arrangement of a correct and throughout data model. To fulfill an analysis, the volume of information should be sorted out according to the commonalities. The main question is, what commonality parameter provides the best results and what is implicated under "the best" definition at all.

Four Basic Algorithms And How To Choose One

1. Clusterization, based on the computation of distances between the objects of the whole dataset, is called connectivity-based, or hierarchical. Depending on the "direction" of the algorithm, it can unite or, inversely, divide the array of information – the names agglomerative and divisive appeared from this exact variation. The most popular and reasonable type is the agglomerative one, where you start by inputting the number of data points, that then are subsequently united into larger and larger clusters, until the limit is reached.

Example - Hierarchical clustering algorithm

2. Centroid-based clustering, is the most frequently occurred model thanks to its comparative simplicity. The model is aimed at classifying each object of the dataset to the particular cluster. The number of clusters (k) is chosen randomly, which is probably the greatest "weakness" of the method. This k-means algorithm is especially popular in machine learning thanks to the likeness with k-nearest neighbours (kNN) method.

Example - K-Means algorithm

3. Expectation-maximization algorithm, at the same time, allows avoiding those complications while providing an even higher level of accuracy. Simply put, it calculates the relation probability of each dataset point to all the clusters we've specified. The main "tool" that is used for this clusterization model is the Gaussian Mixture Models (GMM) – the assumption that the points of the dataset generally follow the Gaussian distribution. The k-means algorithm is, basically, a simplified version of the EM principle. They both require manual input of clusters number, and that's the main intricacy the methods bear. Apart from that, the principles of computing (either for GMM or k-means) are simple: the approximate range of cluster is specified gradually with each new iteration.

Example - K-Means algorithm

4. Finally, density-based clustering comes. The name comprises the main point of the model – to divide the dataset into clusters the counter inputs the parameter, the "neighbourhood" distance. If the object is located within the circle (sphere) of the radius, it, therefore, relates to the cluster.

Example - DBSCAN algorithm

Measures for comparing clustering algorithms

1. Internal measures, which uses intrinsic information in the data to assess the quality of the clustering. Internal measures include the connectivity, the silhouette coefficient and the Dunn index
2. Stability measures, a special version of internal measures, which evaluates the consistency of a clustering result by comparing it with the clusters obtained after each column is removed, one at a time.