

## Problem assignment 9

*Due: Thursday, April 8, 2021*

### Problem 1. K-means clustering

Please load the dataset *clustering\_data.txt*.

**Part a.** Run the k-means algorithm (implemented in Matlab in the function `kmeans`) for finding 3 clusters. Use Euclidean distance to define the differences in between the points. Report the sizes of the three groups found by the `kmeans`. Use scatter function to plot the data in the dataset and the means of the clusters. Please use colors to distinguish data that were assigned different groups. Use a separate color to show the cluster centers (means). Include the plot in your report.

**Part b.** Repeat the setup in Part a. but now assume the number of means is 4. Again report the sizes of the groups and plot the data (with different group colors) and the means found by `kmeans`.

**Part c.** The `kmeans` procedure (if initial means seeds are not set) uses a random set of seeds in each run. Rerun `kmeans` algorithm for  $k = 4$  (the same as Part b). The means found are likely to change. If they did not, try to rerun the procedure again till you see the change in the means. Show the scatter plot of the results when the centers changed.

**Part d.** Let us assume the two runs of the k-means lead to two different clusterings. Write a math expression that would let you compare these different clusterings and pick the best one. Hint: what criterion does the k-means optimize?

**Part e.** Run the `kmeans` procedure (in the default mode) with  $k = 4$  30 times. Report the cluster sizes found for these different runs? Use formula from Part d to decide which clustering is the best. Show the scatter plot of the best clustering.

### Problem 2. Hierarchical clustering

Please load the dataset *clustering\_data.txt* and keep it in variable  $Y$ .

Part a. Run matlab's `linkage` function to create a hierarchy of clusters (cluster tree):  
 $Z = \text{linkage}(Y, 'complete', 'euclidean');$

This will create a hierarchy of clusters using the euclidean distance for pairs of points and 'max' distance for linkages. Plot the dendrogram of the full cluster tree using function `dendrogram(Z,0)`; Include the graph in the report.

Part b. The cluster tree can be used to define clusterings with the different number of groups. Use function `cluster`:

```
C = cluster(Z,'maxclust',4);
```

to assign data instances to four clusters. Using scatter function plot the results obtained by the hierarchical clustering. Similarly to Problem 2 use colors to distinguish the groups found. Include the graph in your report. Compare the results to Problem 2. part e. Are the clusters the same or different?

### Problem 3. Feature/Input ranking

Consider the dataset in file *FeatureSelectionData.txt*. The dataset consists of 259 examples (in rows) where each example is defined by 70 dimensional input vector (represented in columns) and an associated binary label (in last column).

**Part a.** Write and submit a function *Fisher\_score(x,y)* that takes as arguments a vector of one-dimensional inputs  $x$  and a vector of binary outputs  $y$  and calculates the Fisher score as defined in the lecture. Use this function to evaluate the different dimensions of the input space (there are 70 dimensions) to estimate their individual predictive power. Please report the ordered list of dimensions with the top 20 Fisher scores, and their Fisher score values. The dimensions should be labeled from 1 to 70 depending on their position in the dataset.

**Part b.** Write and submit a function *AUROC\_score(x,y)* that takes as arguments a one-dimensional vector of inputs  $x$  and a vector of outputs  $y$  and calculates the area under the ROC curve. You may use Matlab functions to calculate the area under the curve for this purpose. Similarly to part a, evaluate the different dimensions of the input space and their individual predictive power based on AUROC score. Again, report the ordered list of 20 dimensions with the top 20 AUROC scores, and their values. Compare the results from part a and part b and discuss your findings. Are the ordered lists the same? In general, do you expect them to be the same.