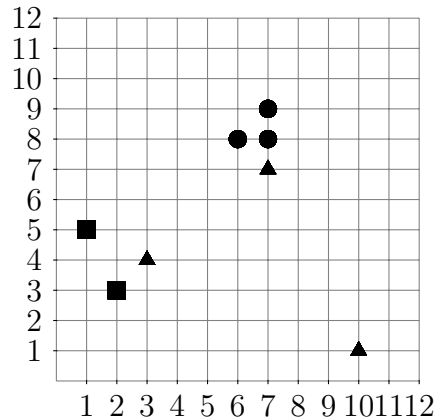


Exercise 5 : Clustering : k -means and Silhouette, Evaluation of Classifiers, Probability

Exercise 5-1 : k -means, choice of k , and compactness

Given the following data set with 8 objects (in \mathbb{R}^2) as in the lecture :



Compute a complete partitioning of the data set into $k = 3$ clusters using the basic k -means algorithm (due to Forgy and Lloyd). The initial assignment of objects to clusters is given using the triangle, square, and circle markers.

Objects x are assigned to the cluster with the least increase in squared deviations $SSQ(x, c)$ where c is the cluster center.

$$SSQ(x, c) = \sum_{i=1}^d |x_i - c_i|^2$$

Start with computing the initial centroids, and draw the cluster assignments after each step and explain the step. Remember to use the least squares assignment !

Give the final quality of the clustering (TD^2). How does it compare with the solutions for $k = 2$ discussed in the lecture ? Can we conclude on $k = 3$ or $k = 2$ being the better parameter choice on this data set ?

Also compute solutions with $k = 4$, $k = 5$, starting from some random initial assignments of objects to clusters. What do you observe in terms of the TD^2 measure ?

Exercise 5-2 : Measure for Evaluation of Classifiers

Given a data set with known class labels ($f(o)$) of the objects. In order to evaluate the quality of a classifier h , each object is additionally classified using h . The results are given in the table (all three columns) below.

ID	$f(o)$	$h(o)$
O_1	A	A
O_2	B	A
O_3	A	C
O_4	C	C
O_5	C	B

ID	$f(o)$	$h(o)$
O_6	B	B
O_7	A	A
O_8	A	A
O_9	A	A
O_{10}	B	C

ID	$f(o)$	$h(o)$
O_{11}	B	A
O_{12}	C	A
O_{13}	C	C
O_{14}	C	C
O_{15}	B	B

- Rewrite the definitions for precision and recall given in the lecture by using TP, TN, FP, and FN.
- Using the table (all three columns) above, compute precision and recall for each class.
- To get a complete measure for the quality of the classification with respect to a single class, the F_1 -measure (the harmonic mean of precision and recall) is commonly used. It is defined as follows :

$$F_1(h, i) = \frac{2 \cdot \text{Recall}(h, i) \cdot \text{Precision}(h, i)}{\text{Recall}(h, i) + \text{Precision}(h, i)}$$

Compute the F_1 -measure for all classes.

- So far, the F_1 -measure is only defined for classes and not yet useful to get an overview of the overall performance of the classifiers. To achieve such an overall assessment, one commonly takes the average over all classes using one of the following two approaches :
 - Micro Average F_1 -Measure : The values of TP , FP and FN are added up over all classes. Then precision, recall and F_1 -measure are computed using these sums.
 - Macro Average F_1 -Measure : Precision and recall are computed for each class individually, afterwards the average precision and average recall are used to compute the F_1 -measure.

Compute the Micro- and Macro-Average F_1 -measures for the example above. What do you observe?

Exercise 5-3 : Procedures for Evaluation of Classifiers

Given a data set D with objects from classes A and B ($D = A \cup B$) where the class assignments are *random* (not related to the attribute values). Furthermore, let the two classes have the same size $|A| = |B|$.

- (a) What *true error rate* is to be expected for an *optimal* (for this data set) classifier?
- (b) What error rates are to be expected when training and evaluating an optimal classifier on the given dataset using a leave-one-out test?
- (c) Remember that in Bootstrap we produce the training and test data by sampling with replacement. An object is with a probability of

$$\left(1 - \frac{1}{n}\right)^n \approx 0.368$$

not part of the n training objects, i.e. only about 63.2% of the objects are used for training. (Compare this to 10-fold cross validation, where 90% of the data are used for training.)

This implies that the error estimation is pessimistic, as the training set has size n , but actually only contains $0.632 \cdot n$ *different* examples.

To make up for this, when evaluating bootstrap it is a common practice to also include the apparent classification error (error on the training data) during evaluation :

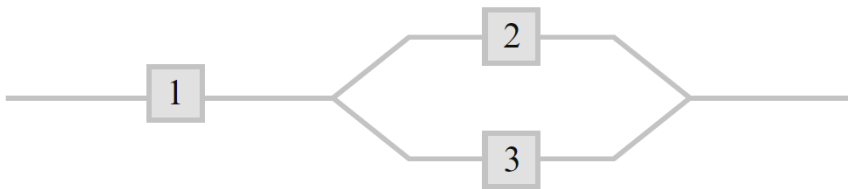
$$\text{error rate} = 0.632 \cdot \text{Error on test set} + 0.368 \cdot \text{Error on training set}$$

This will be repeated multiple times (with different samples) and averaged.

What error rates are to be expected when evaluating an optimal classifier on the given dataset using the 0.632 Bootstrap method? Interpret these results.

Exercise 5-4 : Events and Sample Spaces

- (a) We have a system of several fuses. We can examine each single fuse to see whether it is defective. The sample space for this experiment can be abbreviated as $\Omega = \{N, D\}$, where N represents not defective, D represents defective.
- If we examine three fuses in sequence and note the result of each examination, what is the sample space Ω ?
- (b) As an experiment, we observe the number of pumps in use at a six-pump gas-station, so simple events are the numbers 0-6 (pumps in use). Given the events $A = \{0, 1, 2, 3, 4\}$, $B = \{3, 4, 5, 6\}$, and $C = \{1, 3, 5\}$, which simple events are contained in
- $A \cup B$?
 - $A \cup C$?
 - $A \cap B$?
 - $A \cap C$?
 - \bar{A} ?
 - $\overline{A \cup C}$?
- (c) Three components are connected to form a system as shown in this diagram :



Because the components in the 2-3 subsystem are connected in parallel, that subsystem will function if at least one of the two individual components functions. For the entire system to function, component 1 must function and so must the 2-3 subsystem. The experiment consists of determining the condition of each component (S (success) for a functioning component and F (failure) for a non-functioning component).

- What outcomes are contained in the event D that exactly two out of the three components function?
- What outcomes are contained in the event E that at least two of the components function?
- What outcomes are contained in the event G that the system functions?
- List the outcomes in \bar{G} , $D \cap G$, $D \cup G$, $E \cup G$, and $E \cap G$.

Exercise 5-5 : Tools : Over-fitting , Under-fitting

- (a) Load python packages : `numpy`, `matplotlib.pyplot`, `datasets`, `metrics` from `sklearn`, `KNeighborsClassifier` from `sklearn.neighbors` and `train-test-split` from `sklearn.model-selection`.
- (b) Load `breast-cancer` dataset and split that into random train and test subsets and assign 40 percent of data to test dataset.
- (c) Train the K Neighbors Classifier model, then predict and plot train and test accuracy for different number of neighbours from 1 to 100.
- (d) Now put the neighbours number equal to 3, and plot train and test accuracy again with 20 percent of data for test size and different train sizes from (0.1 to 0.8).
- (e) Describe which "k" and "train size" for the model can cause over-fitting or under-fitting.