

Data Mining Exercise 2: 03.02.2022

Remember to enroll yourself on the examination of the current course no later than 8 days before an examination date!

1. Let us have a look at file `inco13par.txt`. What are the types of variables and what kind of statistic can be calculated for each variable? How many different diagnoses the data contains? What is the mean age of all patients? You may need function `nanmean()`.
2. Replace all missing variable values (NaN) using respective mean values from each diagnosis class. Diagnose 1 for instance, can be considered to form a cluster in variable space. `Find()` function is useful in this task.
3. Boxplot and histogram are valuable visual tools in evaluation of variable quality and distribution. Use `boxplot()` function and visualize what kind of values variables UVA, US, CYM and PTR have. What can you say about age distribution over all cases in the data? (`hist()`).
4. You can consider the patients as points in given variable space. Calculate Euclidean distances between cases that are in rows 2, 269 and 393. What cases are closest to each other? What kind of problems you may encounter when you use Euclidean distance measure? In addition to Euclidean distance, there are plenty of others. What other distance measures you have heard of?
5. Euclidean distance, as well as other distance measures can be used when we compare similarity between observations. Study the content in file `dm2.m` and calculate differences between pixel values of handwritten numbers that are in positions (5,5), (1,5) and (1,6), where first number is the row index and last number is the column index. Sum up the absolute values of differences. This measure can be considered as distance between two images. What is the main factor that affects the result you got?