# DATA.ML.200 Pattern Recognition and Machine Learning

TAU Signal Processing
*Exercise Week 1: Machine learning fundamentals*

This exercise you need to do without any help from anyone. The exercise verifies that your programming skills and machine learning background are sufficient for the course. This exercise is **mandatory** to pass the course.

1. `python`  Nearest neighbor classifier (20 points)

   a) Write code to load and vectorize the 20 Newsgroups dataset (available in the lecture notebook)

   b) Classify the test samples using a *simple baseline* method for classification. You may implement the baseline yourself or use the Sklearn `DummyClassifier` module using the 'most frequent class' strategy.

   c) Add to your code functionality to measure the computation time of classifying all test samples. Suitable functions can be found from the Python 'timeit' library (for example: *.default_timer()*).

   d) Print the classification accuracy and computation time for the baseline.

   e) Implement a straightforward nearest neigbor classifier that classifies each test sample one by one. At first, Euclidean distances are computed to all training samples, and then, class label of the smallest distance sample is used as a prediction.

   f) It may turn out that your own NN classifier is slow. In that case, use only a portion of the test samples, and use them to estimate the total computation time.

   Print the number of test samples used, estimated total computation time, and the classification accuracy.

   g) Replace your NN classifier with the Sklearn *NearestNeighbor* classifier.

   h) Print the computation time and classification accuracy for the Sklearn NN.

   Return the following items:

   - Python code: onenn.py that does all above.
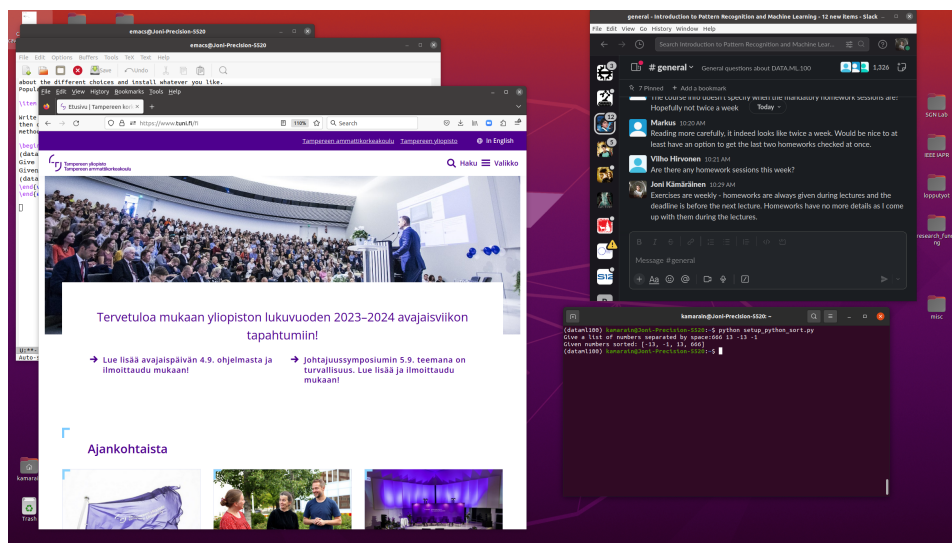   - PNG image: your full desktop screenshot that includes a terminal where the python file is executed and it prints all above numbers: onenn_screenshot.png (see the example)

Figure 1: Example of the required screenshot