### Problem 8.1 (Decision tree)                                          (20 P.)

a) What are the challenges when trying to learn a decision tree from continuous data? (2 P.)

b) Explain why we did not test the same attribute more than once along one path in a decision tree? Under what circumstances and conditions would this make sense? (2 P.)

c) An attribute splits the set of examples $E$ into subsets $E_i$ each having $p_i$ positive and $n_i$ negative examples. Show that this attribute has zero gain if the ratio $\frac{p_i}{(p_i+n_i)}$ is the same for all $i$. Start by defining the gain! (4 P.)

d) Download the "decisiontree.py" file from StudIP and write a decision tree learning algorithm. Complete the stumbs in the file where marked. Train the classifier with the car data set which can be downloaded at https://archive.ics.uci.edu/ml/datasets/Car+Evaluation. (12 P.)

   i) Use 10-fold cross validation and **plot the average accuracies** over the depth of the tree.

   ii) **Illustrate** one learned classifier for depth 3.


### Problem 8.2 (Support vector machines)                                (20 P.)
**Note:** For the SVM we recommend using scikit learn http://scikit-learn.org/stable/.
**Note:** Some subtasks contain multiple questions, give an answer to all!

a) You have a data set with two classes $y_1$ and $y_2$ which differ significantly (factor 1000) in size. **Discuss** the possible consequence for applying an SVM classifier on such data! (4 P.)

b) **State briefly** the basic working of an SVM classifier? What are the critical tuning parameters of SVM classifier. **Explain** their meaning! (4 P.)

c) Explain the fundamentals of the kernel trick. When can it be applied? What benefit does it provide for the application of SVMs? (3 P.)

d) Is the decision boundary affected if the kernel is changed? Explain in either case!(2 P.)

e) Apply the support vector machine classifier (e.g. `sklearn.svm.SVC`) to the IRIS data set. Tune the parameters and **report the corresponding accuracy** of your classifier. What is the best parameter set you found? (4 P.)

f) Compare this result with the other classifiers you have implemented. (3 P.)

---

On the hand-in date, **11.01.2017**, you must hand-in the following: [1]

a) a text file stating how much time you (all together) used to complete this exercise sheet
b) your solutions / answers / code

for problem 8.1 and 8.2.

---

[1] upload via StudIP (if there are problems with the upload contact me **beforehand**: krell@uni-bremen.de)