

Exercise 4

Deadline: 17.12.2020, 16:00.

Ask questions to `#ask-your-tutor-luisa`

This exercise will introduce you to tree-based machine learning methods. Use the stubs provided in file `tree-methods-stubs.ipynb` on MaMPF as a starting point for this homework.

Regulations

Solutions for this week's tasks shall be handed in as a Jupyter notebook `tree-methods-stubs.ipynb`, accompanied with exported HTML `tree-methods-stubs.html`. Zip all files into a single archive `ex04.zip` and upload this file to your assigned tutor on MaMPF before the given deadline.

Note: Each team creates only a single upload, and all team members must *join* it as described in the MaMPF documentation at <https://mampf.blog/zettelabgaben-fur-studierende/>.

Important: Make sure that your MaMPF name is the same as your name on Muesli. We now identify submissions purely from the MaMPF name. If we are unable to identify your submission you will not receive points for the exercise!

1 Density Tree and Decision Tree (20 points)

Complete the code of the classes `DensityTree` and `DecisionTree` in the stub, following the explanations given in the lecture and in code comments. Note in particular, that the prediction of a density should be zero for test points that are located outside the tree's training bounding box.

The optimal splits of a *density tree* shall minimize the leave-one-out error of the resulting children. For a given node l , the leave-one-out error is defined as

$$\text{looErr}_l = \frac{N_l}{NV_l} \left(\frac{N_l}{N} - 2 \frac{N_l - 1}{N - 1} \right)$$

where N_l and V_l are the number of instances in node l and its volume, and N is the total number of instances for the class under consideration. The region spanned by node l is represented by the vectors m_l and M_l of lower and upper bounds respectively, so that $V_l = \prod_{j=1}^D (M_{lj} - m_{lj})$.

In case of a *decision tree*, the optimal splits shall minimize the Gini impurity of the resulting children. The Gini impurity of node l is defined as

$$\text{Gini}_l = N_l \left(1 - \sum_{k=1}^C \frac{N_{lk}^2}{N_l^2} \right)$$

where N_l is the total number of instances in node l , and N_{lk} are the number of instances of class k in l .

To save computation time, only a random subset of size $D_{\text{try}} = \sqrt{D}$ of the features shall be considered when searching for the optimal split in each node. Note that different subsets must be selected in every node. Function `numpy.random.permutation()` and “advanced indexing” (see <https://docs.scipy.org/doc/numpy/reference/arrays.indexing.html#advanced-indexing>) will be helpful here.

Candidate thresholds shall be placed in the middle between consecutive feature values

$$t_{i'j} = \frac{X_{[i]j} + X_{[i+1]j}}{2}$$

provided that $X_{[i]j}$ and $X_{[i+1]j}$ are different. The brackets in $X_{[i]j}$ denote sorted order w.r.t. to feature j . If the two feature values are the same, no threshold shall be placed there. Describe in a comment why this is necessary.

2 Evaluation of Density Tree and Decision Tree (6 points)

Once again, use the digits dataset provided by sklearn. This time, we will not split the data into training and test sets and just measure the training error (otherwise, the dataset would become too small for density estimation). Train a generative classifier using one instance of `DensityTree` per digit class, and a discriminative classifier using one instance of `DecisionTree`. You may try different values for the hyperparameter `n_min` to improve performance. For each method, show the full training error confusion matrix and comment on the results.

3 Density Forest and Decision Forest (8 points)

Complete the code of the classes `DensityForest` and `DecisionForest` in the stub, following the explanations given in the lecture and in code comments. To make the trees in the forest independent of each other, create a new bootstrap training set for each tree. That is, each new training set has the same size as the original one and is obtained by selecting instances from the original training data uniformly at random *with replacement*. Functions `numpy.random.choice()` and “advanced indexing” are useful here.

In case of the `DecisionForest`, it makes sense to drop the size restriction of the split nodes, i.e. to train all leaves to purity, by setting hyperparameter `n_min = 0`.

The ensemble prediction shall be the average of the individual tree predictions.

4 Evaluation of Density Forest and Decision Forest (6 points)

Train a generative classifier using 10 instances of `DensityForest` and a discriminative classifier using one instance of `DecisionForest`, with each forest consisting of 20 trees. For each method, show the full training error confusion matrix and comment on the results.

Compare your results with the confusion matrix obtained from sklearn’s predefined decision forest `sklearn.ensemble.RandomForestClassifier`, also trained with 20 trees.