I this lab I learned decision tree classifiers and performed 5-fold cross-validation on two different datasets using Python's sklearn library.

**<u>Data sets description and missing values:</u>**
<u>congressional voting records data set</u>
The first data set was a congressional voting records data set, which contains votes of republican and democrat U.S. House Representatives on 16 key issues. Therefore, there are two classes: republican and democrat and 16 attributes that take on binary values yes (1) or no (0). There are also some missing data in this dataset. In particular, there are 16 observations (rows) that contain some missing values. There are various options how to handle missing data. Wen can either drop the observations that contain missing values from the dataset, or we can impute the values. In the case of the congressional voting dataset, there are 435 observations in total and only 16 (3.7%) of them contain missing data. Since the number of rows with missing data is very small compared to the total dataset I decided to simply drop those rows from the dataset.
<u>breast cancer data set</u>
The second dataset was a breast cancer data set, which contains numerical attributes and two classes: malignant and benign. The dataset has 699 rows in total out of which 203 rows contain missing data. That represents 29%, which is a lot and dropping them would significantly reduce the amount of training data, which in turn could have a negative effect on performance of the classifier. This is why I decided to impute the missing data using sklearn's SimpleImputer and choosing a the "most_frequent" strategy.

**<u>5-fold cross-validation and 95% confidence interval:</u>**
5-fold cross-validation divides a dataset into 5 datasets of equal size and applies the learning algorithm 5 times. From the 5 datasets each time the learning algorithm is applied, a different dataset is held out as a test data set and the other 4 as the training dataset.

The 95% confidence interval was computed using the z score and standard deviation as according to (Mitchell, 1997).
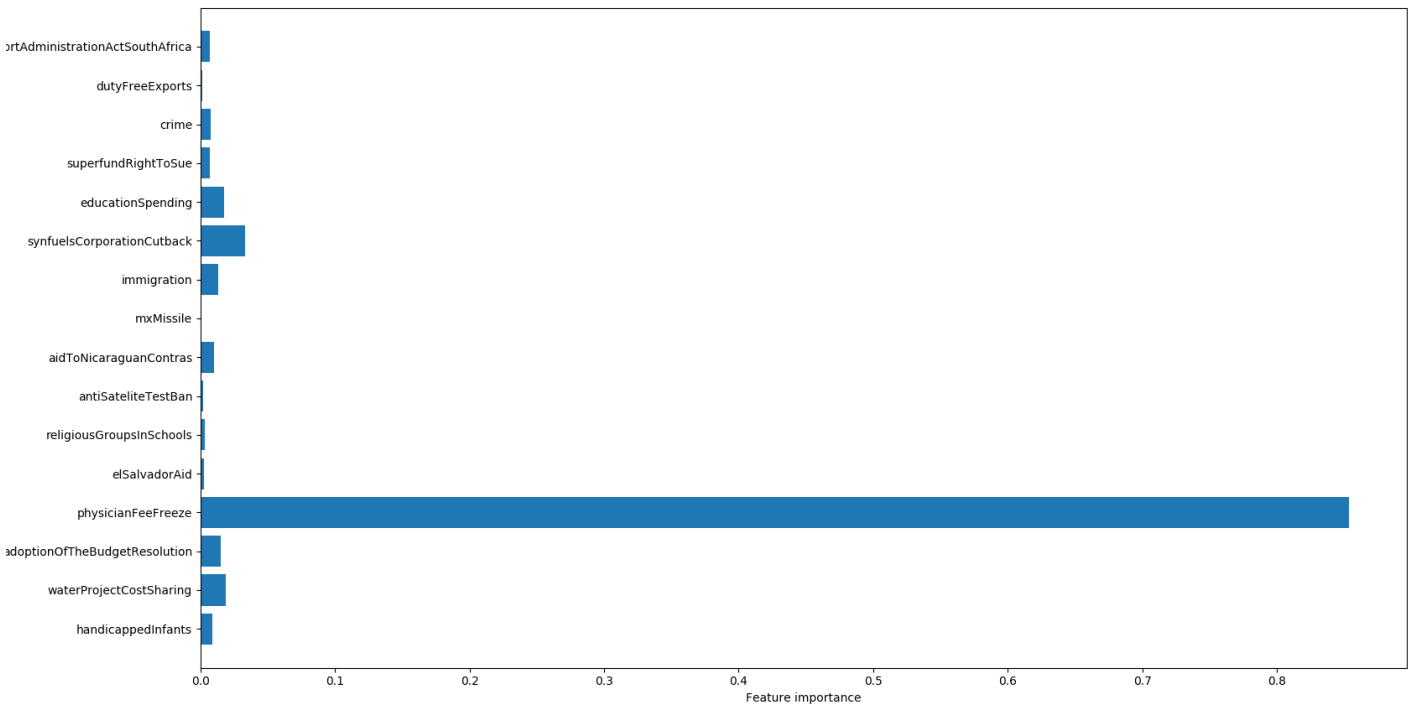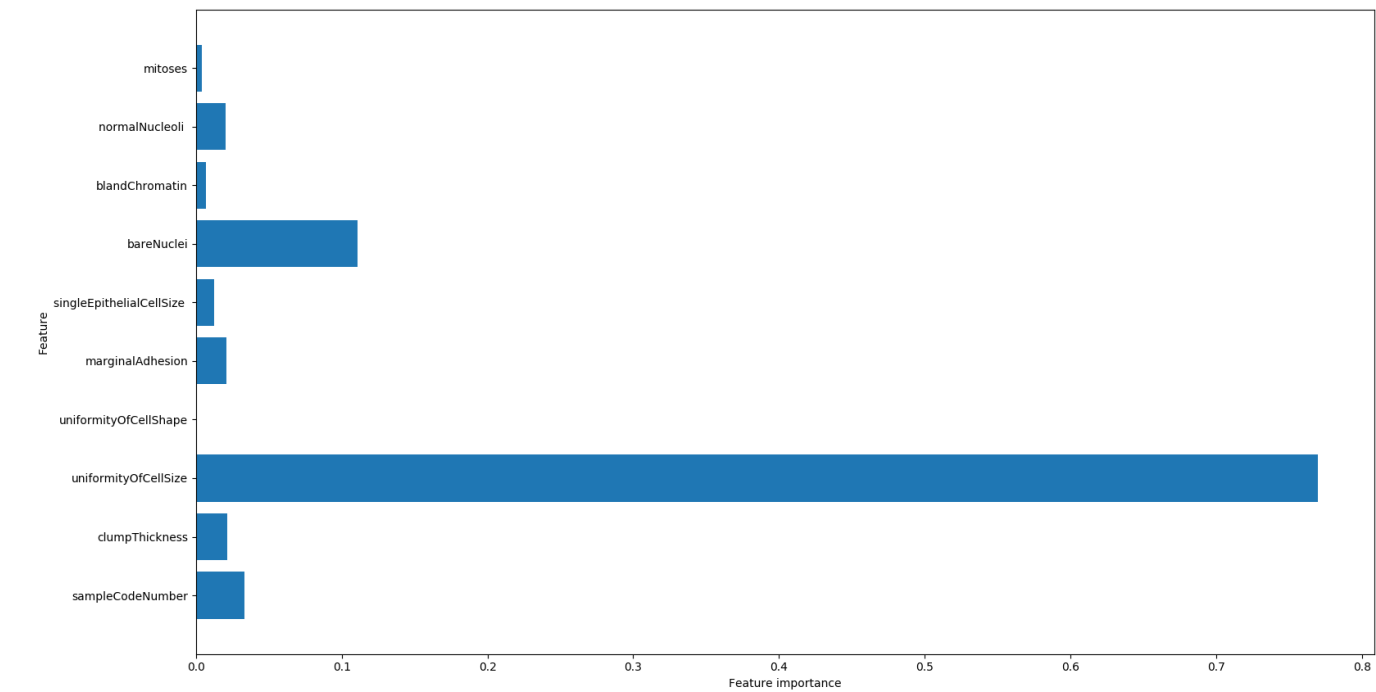z at 95% confidence interval = 1.96

standard deviation of accuracy (h) = $\sqrt{\dfrac{accuracy\,(h)*(1-accuaracy\,(h))}{n}}$

95% confidence interval = accuracy(h) +/- (z * standard deviation)

**<u>TASK 1:</u>**
<u>Interesting Rules:</u>
I visualized the importance of each feature on a scale of 0 to 1 (0 = feature not important at all, 1 = feature perfectly predicts the target). In case of the cancer dataset, the most important feature (almost 0.8) is uniformity of cell size. Therefore, it makes sense that this attribute is the root node of the cancer decision tree classifier. The second most important feature is bare nuclei (around 0.1), which is why it makes sense that this attribute appears in the first level of the tree right below the root node. The most important feature (over 0.8) of the voting dataset is physician fee freeze. Therefore, it makes sense that this feature is the root node of the voting decision tree classifier since it provides the most amount of information. The other features have little importance. Therefore, I expect that the five trees in task 3 (c) will have the same root node, but may differ in the lower levels.

physicianFeeFreeze <= 0.5
samples = 304
value = [178, 126]
class = republican

True / False

superfundRightToSue <= 0.5
samples = 170
value = [168, 2]
class = republican

synfuelsCorporationCutback <= 0.5
samples = 134
value = [10, 124]
class = democrat

samples = 122
value = [122, 0]
class = republican

synfuelsCorporationCutback <= 0.5
samples = 48
value = [46, 2]
class = republican

adoptionOfTheBudgetResolution <= 0.5
samples = 108
value = [1, 107]
class = democrat

educationSpending <= 0.5
samples = 26
value = [9, 17]
class = democrat

religiousGroupsInSchools <= 0.5
samples = 27
value = [25, 2]
class = republican

samples = 21
value = [21, 0]
class = republican

samples = 93
value = [0, 93]
class = democrat

immigration <= 0.5
samples = 15
value = [1, 14]
class = democrat

immigration <= 0.5
samples = 7
value = [5, 2]
class = republican

adoptionOfTheBudgetResolution <= 0.5
samples = 19
value = [4, 15]
class = democrat

samples = 12
value = [12, 0]
class = republican

waterProjectCostSharing <= 0.5
samples = 15
value = [13, 2]
class = republican

aidToNicaraguanContras <= 0.5
samples = 4
value = [1, 3]
class = democrat

samples = 11
value = [0, 11]
class = democrat

samples = 4
value = [4, 0]
class = republican

handicappedInfants <= 0.5
samples = 3
value = [1, 2]
class = democrat

superfundRightToSue <= 0.5
samples = 14
value = [1, 13]
class = democrat

waterProjectCostSharing <= 0.5
samples = 5
value = [3, 2]
class = republican

samples = 6
value = [6, 0]
class = republican

dutyFreeExports <= 0.5
samples = 9
value = [7, 2]
class = republican

samples = 1
value = [1, 0]
class = republican

samples = 3
value = [0, 3]
class = democrat

samples = 1
value = [1, 0]
class = republican

samples = 2
value = [0, 2]
class = democrat

immigration <= 0.5
samples = 2
value = [1, 1]
class = republican

samples = 12
value = [0, 12]
class = democrat

samples = 2
value = [0, 2]
class = democrat

samples = 3
value = [3, 0]
class = republican

crime <= 0.5
samples = 6
value = [4, 2]
class = republican

samples = 3
value = [3, 0]
class = republican

samples = 1
value = [1, 0]
class = republican

samples = 1
value = [0, 1]
class = democrat

samples = 1
value = [0, 1]
class = democrat

antiSateliteTestBan <= 0.5
samples = 5
value = [4, 1]
class = republican

samples = 2
value = [2, 0]
class = republican

elSalvadorAid <= 0.5
samples = 3
value = [2, 1]
class = republican

samples = 1
value = [1, 0]
class = republican

samples = 2
value = [1, 1]
class = republican

**TASK 2:**

Breast cancer data set accuracy
Accuracy based on 95 percent confidence interval: average 5-fold cross validation score:
0.9326953230602866 +/- 0.01879049587007505

5-fold cross validation scores:
        [0.91970803 0.90510949 0.9270073  0.96350365 0.94814815]

Simple accuracy (not k–fold cross-validation estimation accuracy):
        Accuracy on training set: 1.0
        Accuracy on testing set: 0.9463414634146341


**TASK 3:**

Congressional data set accuracy
Accuracy based on 95 percent confidence interval: average accuracy on testing 0.9241379310344827
+/- 0.02488239513283033

Accuracies for $D_i$ for i=1, 2, ...5:
        k = 1
        Accuracy on training set: 0.9942528735632183
        Accuracy on testing set: 0.9195402298850575

        k = 2
        Accuracy on training set: 0.9971264367816092
        Accuracy on testing set: 0.9655172413793104

        k = 3
        Accuracy on training set: 0.9971264367816092
        Accuracy on testing set: 0.9080459770114943

        k = 4
        Accuracy on training set: 0.9971264367816092
        Accuracy on testing set: 0.9080459770114943

        k = 5
        Accuracy on training set: 0.9942528735632183
        Accuracy on testing set: 0.9195402298850575

Simple accuracy (not k–fold cross-validation estimation accuracy):
Accuracy on training set: 0.9967105263157895
Accuracy on testing set: 0.9312977099236641

Difference between the trees
The five trees constructed in task 3 and the one constructed in task 1 have the same root node
(physician fee freeze), which makes sense because it is the most important feature by far as discussed
above. The first levels of the trees are fairly similar with adoption of the budget resolution and synfuels

corporation cutback appearing in the first levels. The lower levels somewhat differ, but that is expected since the most important feature by far is physician fee freeze and the others are fairly equally unimportant.