

CourseWork2 Decision Trees

Team MCMC

(Junjie Lin, Soomin Myung, Wenhao Deng, Yinbin Li, Yunfei Wang)

1. Introduction

The purpose of this coursework is to implement a classification and regression tree learning algorithm using MATLAB (basic MATLAB functions). The two datasets trained and tested on the ‘UCI Machine Learning Repository’ are:

- a. Airfoil self-noise data set for regression (Airfoil data),
- b. Estimation of obesity levels based on eating habits and physical condition data set for classification (Obesity data).

Information Gain (entropy) as the splitting criterion was used in the classification tree and Sum-of-Square Residuals (variance) as the splitting criterion for the regression tree.

2. Data Pre-Handling

The data handling part of codes were written in **DataHandling.m**. It performs several functions as listed:

- Download the two dataset files from UCI to the current directory.
- For the classification dataset (Obesity data).
 - a. Dataset has 2111 examples in total, 1500 examples were used for training, and 611 examples were used for testing.
 - b. Transform the multi-class classification to the binary classification. There is a 7-classes label, which indicates the obesity levels of an individual. Based on the coursework requirement, these 7 classes were recategorised to 2 classes (no_obesity and obesity).
 - c. One-hot encode the multi-class features.

Based on the requirement, a binary tree should be built after learning. That means all decision should be made binary (positive or negative) for internal nodes. However, in obesity datasets, there are many multi-class features. Therefore, as shown in figure 1, the discrete value of the original data was encoded to one-hot keys to simplify the training algorithm.

d. Randomise the dataset.

Randomise was done to the data to ensure the distribution of the segment data can reflect the original distribution of the dataset when cross-validation was conducted.

ObesityDataSetRawanddatainthetic																								
Gender	Age	Height	Weight	family_hi...	FAVC	FCVC	NCP	CAEC	SMOKE	CH2O	SCC	FAF	TUE	CALC	MTRANS	NObeyes...								
	▼数值	▼数值	▼数值	▼分类	▼分类	▼数值	▼分类	▼分类	▼数值	▼分类	▼数值	▼分类	▼数值	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类
Female	21	1.62	64	yes	no	2	3	Sometimes	no	2	no	0	1	no	Public_Tra...	Normal_W...								
Female	21	1.52	56	yes	no	3	3	Sometimes	yes	3	yes	3	0	Sometimes	Public_Tra...	Normal_W...								
Male	23	1.8	77	yes	no	2	3	Sometimes	no	2	no	2	1	Frequently	Public_Tra...	Normal_W...								
Male	27	1.8	87	no	no	3	3	Sometimes	no	2	no	2	0	Frequently	Walking	Overweig...								
Male	22	1.78	89.8	no	no	2	1	Sometimes	no	2	no	0	0	Sometimes	Public_Tra...	Overweig...								
Male	29	1.62	53	no	yes	2	3	Sometimes	no	2	no	0	0	Sometimes	Automobile	Normal_W...								
Female	23	1.5	55	yes	yes	3	3	Sometimes	no	2	no	1	0	Sometimes	Motorbike	Normal_W...								
Male	22	1.64	53	no	no	2	3	Sometimes	no	2	no	3	0	Sometimes	Public_Tra...	Normal_W...								

ProcessedDailyDataset																									
FAF	TUE	Gender_Fe...	Gender_M...	family_hi...	family_hi...	FAVC_no	FAVC_yes	CAEC_Alw...	CAEC_Freq...	CAEC_So...	CAEC_no	SMOKE_no	SMOKE_yes	SCC_no	SCC_yes	CALC_Alw...	CALC_Freq...	CALC_So...	CALC_no	MTRANS_...	MTRANS_...	MTRANS_...	MTRANS_...	MTRANS_...	NOkeyes...
▼数值	▼数值	▼分类	▼分类	▼分类	▼分类	▼数值	▼数值	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类	▼分类
2	1	0	1	0	1	1	0	0	0	1	0	1	0	1	0	1	0	0	0	0	0	1	0	0	
1	1	0	0	1	0	1	0	0	1	0	1	0	1	0	0	0	1	0	0	0	0	1	0	0	
0.880712	0	1	0	0	1	0	1	0	0	1	0	1	0	1	0	0	0	1	0	0	0	0	1	0	
0.174475	0.261705	0	1	0	1	0	1	0	0	1	0	1	0	1	0	0	0	0	1	0	0	0	1	0	
0	1	0	1	0	1	0	1	0	0	1	0	1	0	1	0	0	0	0	1	0	0	0	1	0	
0	2	0	1	0	1	0	1	0	0	1	0	1	0	1	0	0	0	1	0	0	0	0	1	0	
0.998391	0.85882	0	1	1	0	0	1	0	0	1	0	1	0	1	0	0	0	1	0	0	0	0	1	0	
0.008127	1	1	0	0	1	0	1	0	0	1	0	1	0	1	0	0	0	1	0	0	0	0	1	0	
1.862235	1	1	0	0	1	0	1	0	0	1	0	1	0	1	0	0	0	1	0	0	0	0	1	0	
2	0.152985	1	0	0	1	0	1	0	1	0	0	1	0	1	0	0	0	0	1	0	0	0	1	0	
0	1.119877	0	1	0	1	0	1	0	0	1	0	1	0	1	0	0	0	0	1	0	0	0	1	1	

Figure 1. Comparison between original data and processed data

- For the regression dataset (Air-foil data).

- a. Dataset has 1503 examples in total, 1000 examples for training, and 502 examples for testing.
- b. Randomise the dataset.

Fortunately, there is no missing value in these two datasets, so how to deal with missing value data was not considered in the coursework. However, if the dataset did have missing value examples and if the dataset was big enough, the missing value examples could be deleted in the data pre-handling process.

Finally, after the handling of the data, the function would return classification features and its labels, also regression features and its labels.

3. Results

The main code file is **DecisionTreeLab.m**, which could be run to see the results and the code file **DecisionTreeLearning.m** is the interface implemented

following the guidance of the lab sheet.

- For the classification

Splitting criterion: Information Gain (entropy).

Leaf node determine condition: All examples have the same label value.

1. The first learning with the training set and test set (Tree is shown in figure 2).

[Training] Precision: 1 Recall: 1 F1-score: 1

[Test] Precision: 0.99303 Recall: 0.98276 F1-score: 0.98787

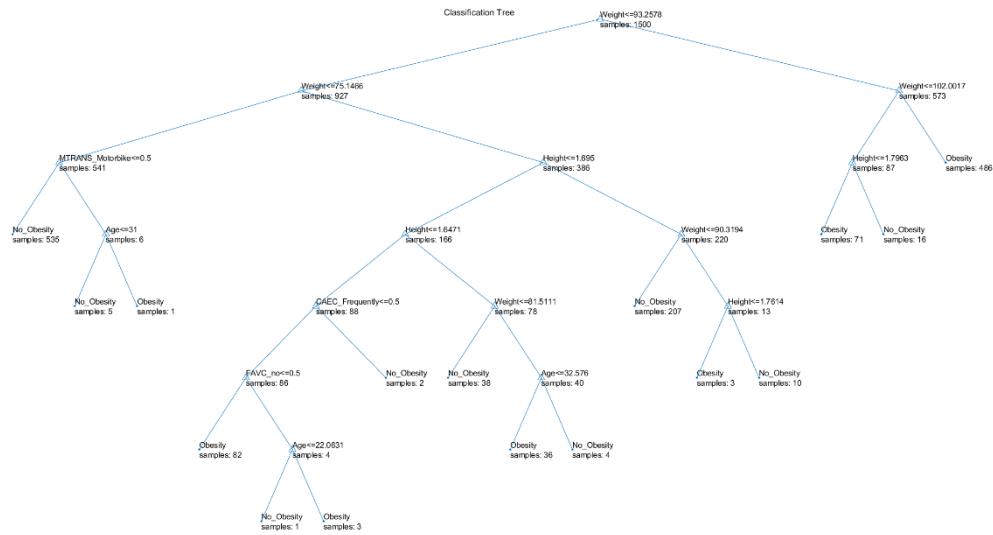


Figure 2. Classification tree learning from the training set.

2. 10-fold cross-validation (Trees are shown in figure 3, also see appendix a).

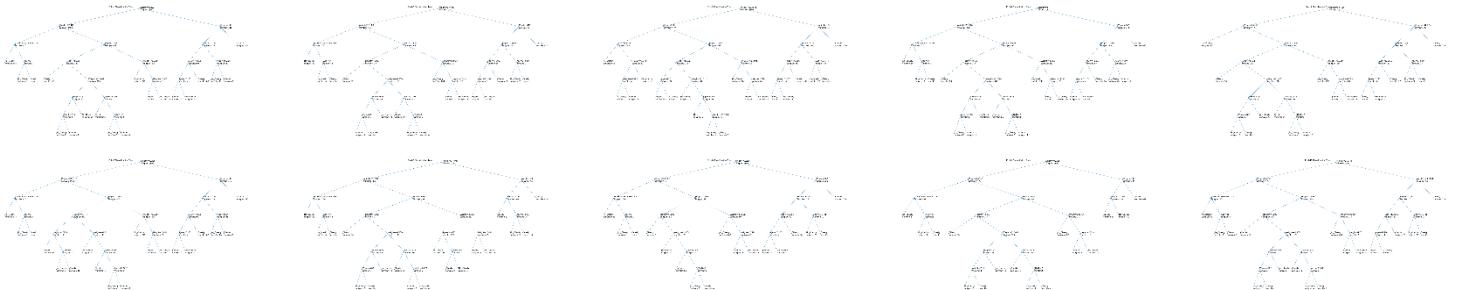


Figure 3. Trees learned from 10-fold cross-validation

[Training] Average Precision: 1 Average Recall: 1 Average F1-score: 1

[Test] Average Precision: 0.99362 Average Recall: 0.98891

Average F1-score: 0.99121

- For the regression

Splitting criterion: Residuals Sum-of-Square (variance).

Leaf nodes determine conditions:

- A. All examples have the same label value.
 - B. The number of examples equals or less than 9.
 - C. The based residuals sum-of-square equals or less than 9.
1. First learning with the training set and test set (Tree is shown in figure 4).
- RMSE of training set: 1.8941
- RMSE of test set: 3.0047

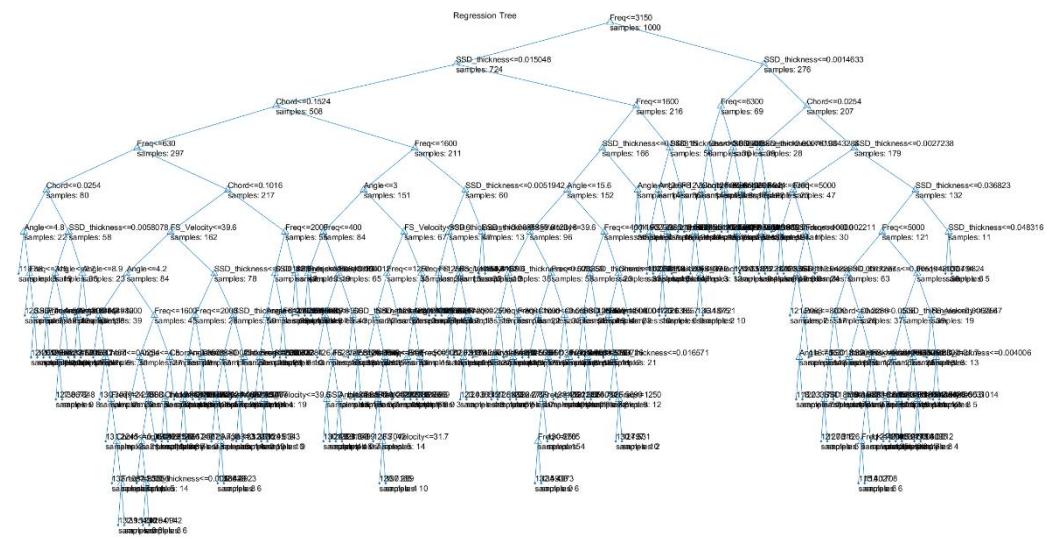


Figure 4. Regression tree learning from the training set.

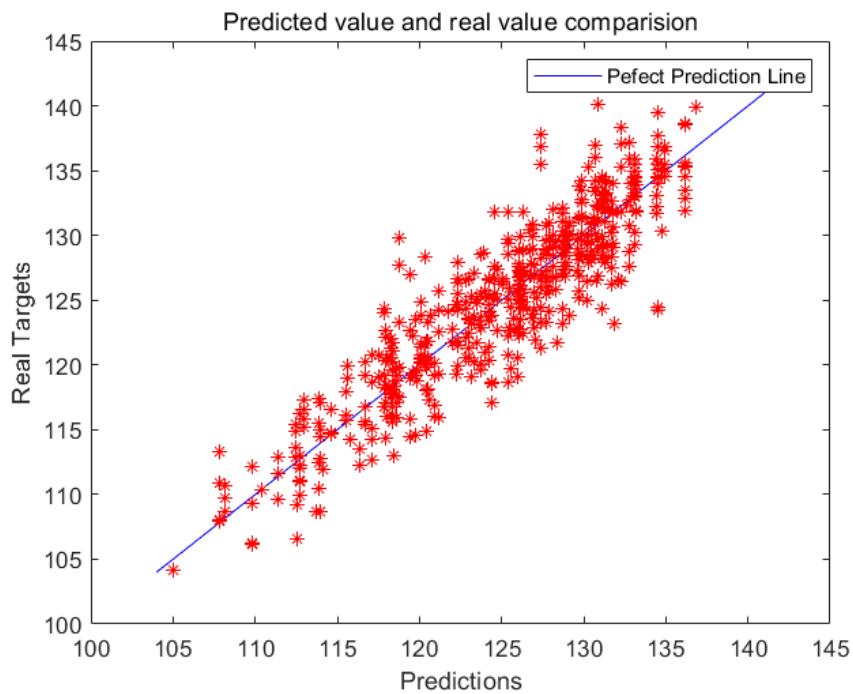


Figure 5. Predicted value vs real targets.

2. 10-fold cross-validation (10-fold Tree drawings see appendix b).

Average RMSE of training sets: 1.8457

Average RMSE of validation sets: 2.9665

3. Experiments with hyperparameters

New leaf nodes determine condition was added after the 10-fold cross-validation: the max depth of the tree.

As shown in figure 6, the RMSE of the testing set went flat following the increasing of tree depth after 11.

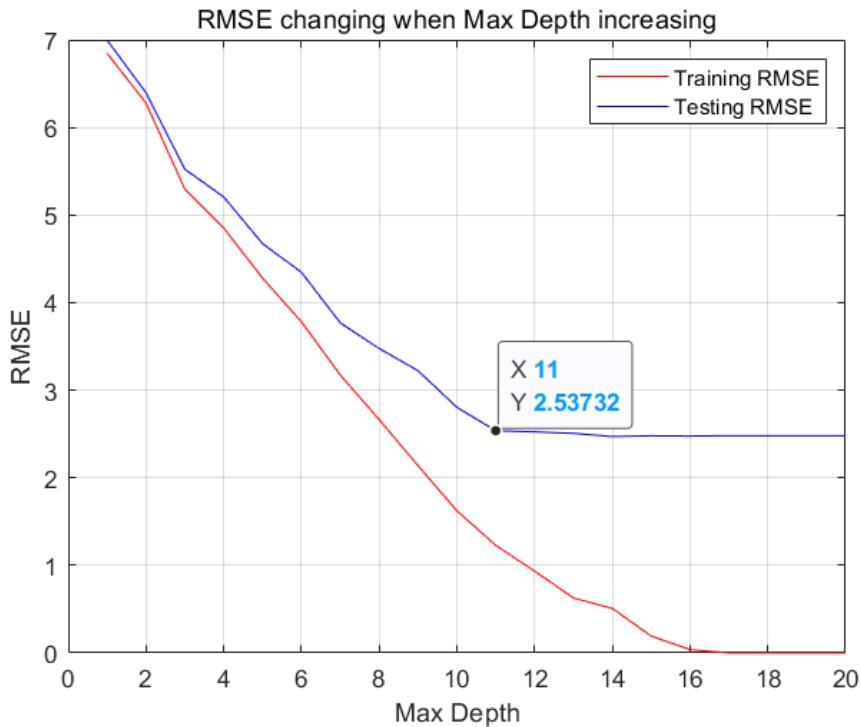


Figure 6. RMSE changing when Max Depth increasing

Therefore, we set the max depth as 11, minimal RSS and minimal numbers of examples were the same as the first training. As shown in figure 7, a tree which is smaller than the previous one was learned, and without losing accuracy which shown in figure 8 (RMSE of training set: 1.9335 RMSE of test set: 3.0084).

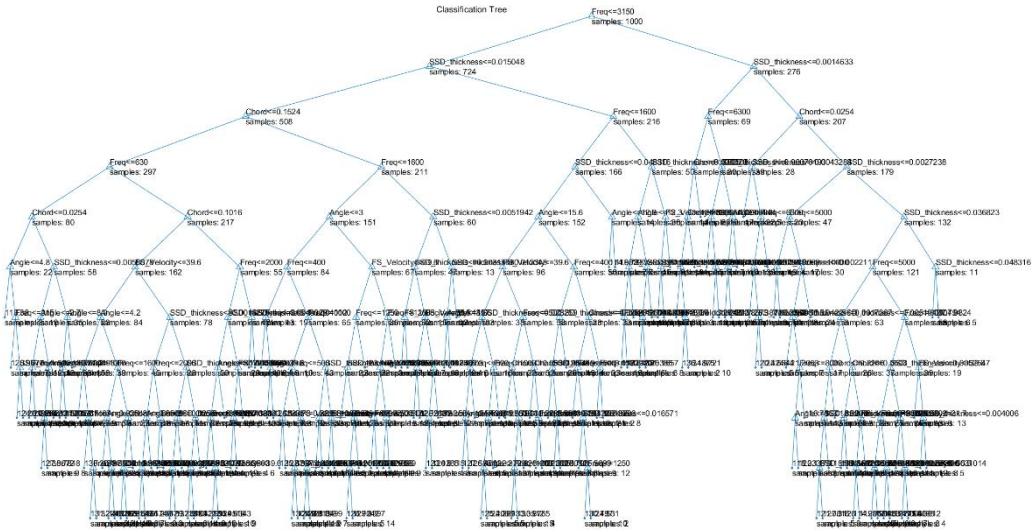


Figure 7. The regression tree with max depth as 11.

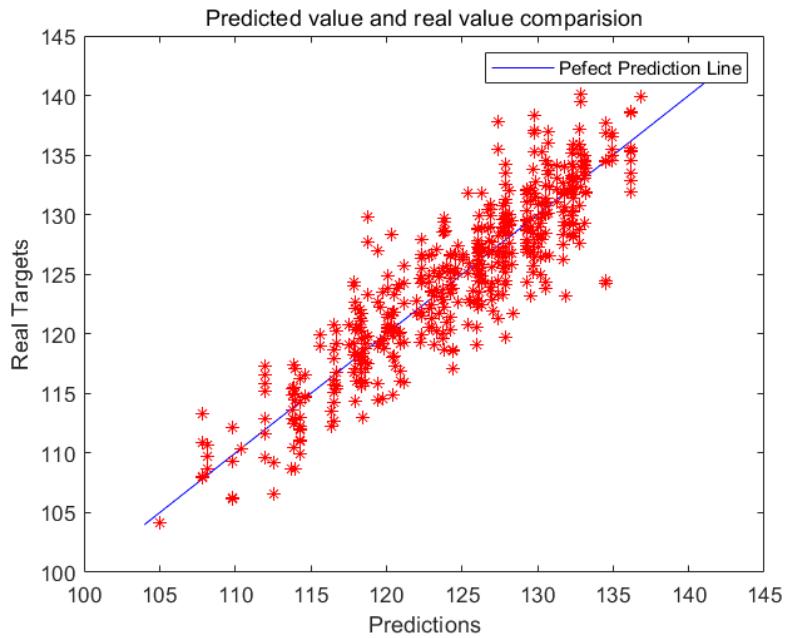


Figure 8. Predicted value vs real targets in max depth as 11.

4. Additional question answers

1. Pruning

Pruning can reduce the size of the decision tree by deleting the redundant subtrees. By reducing the complexity of the tree, pruning can improve the generalisation ability of the tree. According to Brunk and Pazzani in 1991, a technique called REP (Reduce Error Pruning) could be applied to prune the decision tree.

This REP technique split the dataset to be one training set and one pruning set. After using the training set to generate a decision tree, the pruning set is used to test and evaluate which subtree could be removed from the decision tree from the bottom to the top.

Take one of our classification trees as an example:

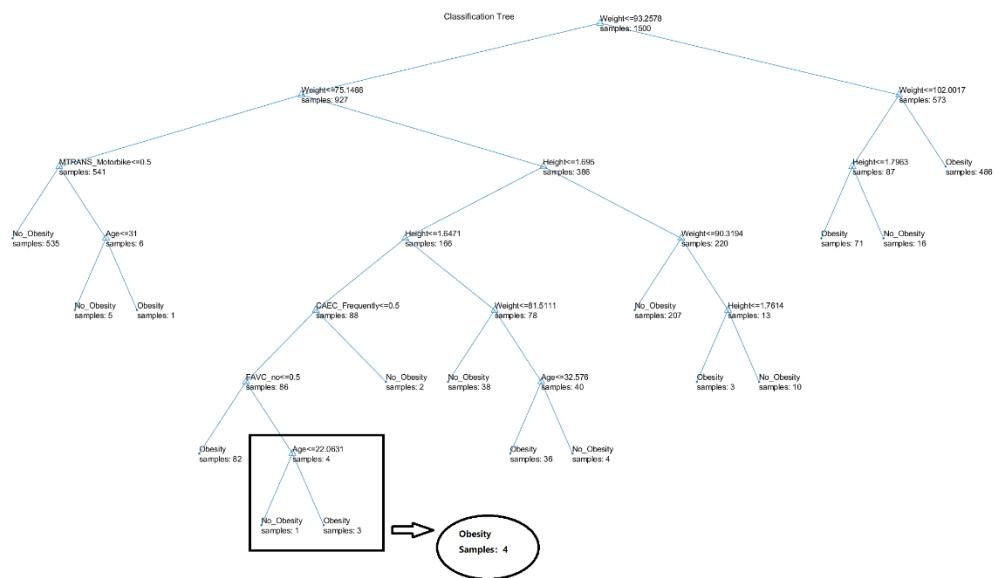


Figure 9. The first step to prune the tree

- Firstly, the pruning set is used to the original tree to produce a based accuracy of the tree.
- After this, the deepest leaf nodes are found and removed. All examples in these two nodes are put into their parent node. As a consequence, the parent node would become a new leaf node, and the prediction would become the majority value of targets corresponding to the examples within the node.

3. Then, this new tree is going to be tested using the pruning set and check whether the accuracy is improved. If the accuracy equals to the based accuracy or it is improved, then the tree stays in the current shape. Otherwise, those two removed leaf nodes should be recovered to the tree.
4. Repletely do step 3 and 4 to all nodes from bottom to the top, and finally, a pruned tree would be generated.

Suppose the original tree has a severe overfitting problem. In that case, the pruned tree should be in a smaller size and would have a higher predictive accuracy or at least has the same accuracy of the original tree.

2. Multiple independent binary labels prediction

A single binary tree could also implement multi-output classification. Although multi-label classification could be transformed to multi-class classification in data pre-handling process, it is not always a good manner to solve multi-label classification. For example, data instance pair $(x_1, [1,1,0])$ could be transformed to pair $(x_1, 6)$ and instance pair $(x_2, [0,1,1])$ could be transformed to pair $(x_2, 3)$, however, the number of transformed classes would grow exponentially when the number of labels grows. Therefore, an adapted algorithm for multi-labels output is worth to be considered.

According to Zhang and Zhou in 2013, an algorithm adaptation method for multi-label classification called Multi-Label Decision Tree was discussed. It also uses information gain as branch splitting criterion. However, to predict multi-binary-labels in a single tree, for all internal nodes in the tree, the entropy formula should be adapted. The entropy computation formula used in this coursework is:

$$Ent = -\frac{p}{p+n} \log_2 \left(\frac{p}{p+n} \right) - \frac{n}{p+n} \log_2 \left(\frac{n}{p+n} \right) \quad (1)$$

For predicting multiple binary labels, because each label is independent with each other, the entropy formula could be computed as the sum of the entropy of each label:

$$MLEnt = \sum_{i=1}^k -\frac{p_i}{p_i+n_i} \log_2 \left(\frac{p_i}{p_i+n_i} \right) - \frac{n_i}{p_i+n_i} \log_2 \left(\frac{n_i}{p_i+n_i} \right) \quad (2)$$

Here, the entropy of k different binary labels is computed one by one and added up together to be multi-labels entropy finally. The rest of the information gain computing is remaining the same as what is done in this coursework.

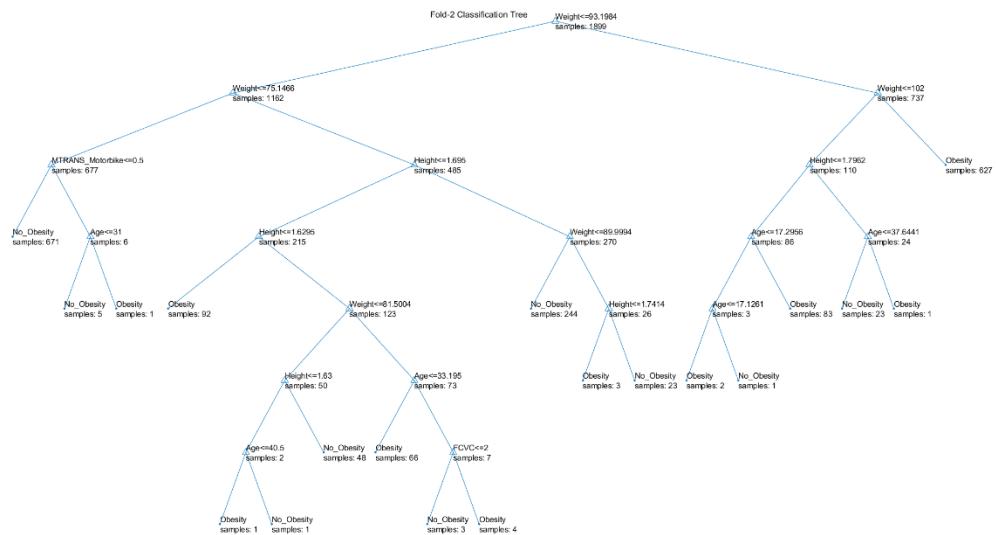
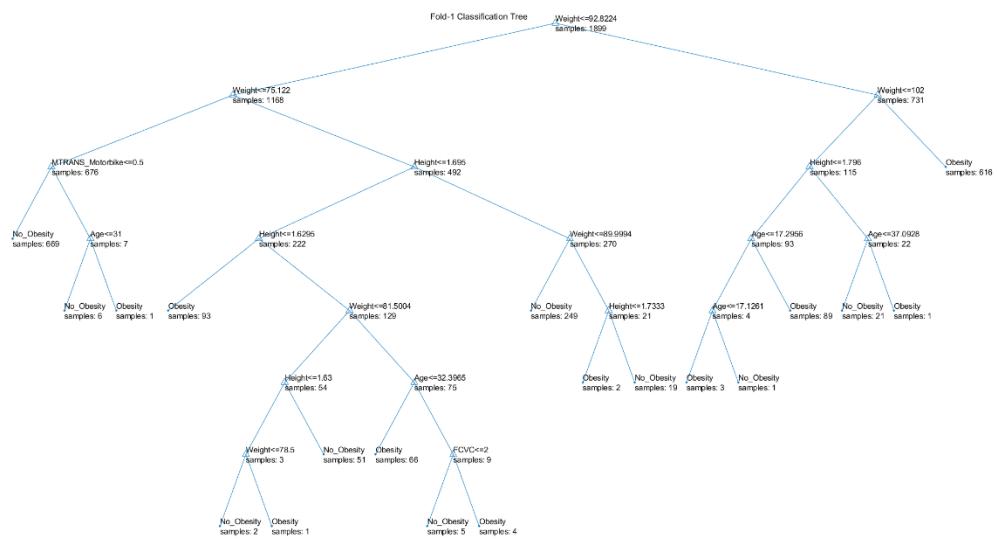
Moreover, what needs to predict would be a set of labels, so when the learning algorithm reaches the leaf node condition, the prediction method would be adapted to use the majority value of each node to compose as a set of labels (prediction). For example, for labels $[l_1, l_2, l_3 \dots l_k]$, first, we find the majority value (positive or negative) of the 1st binary label using all examples within this node, and then find the majority value for the rest of labels one by one. Finally, the decision (prediction) can be made by composing all these majority values found as a set of label values.

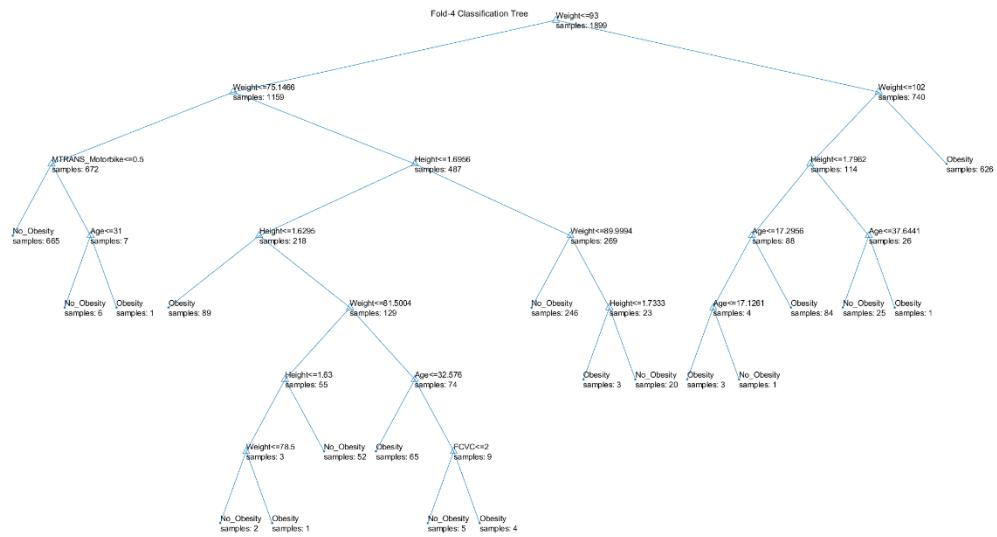
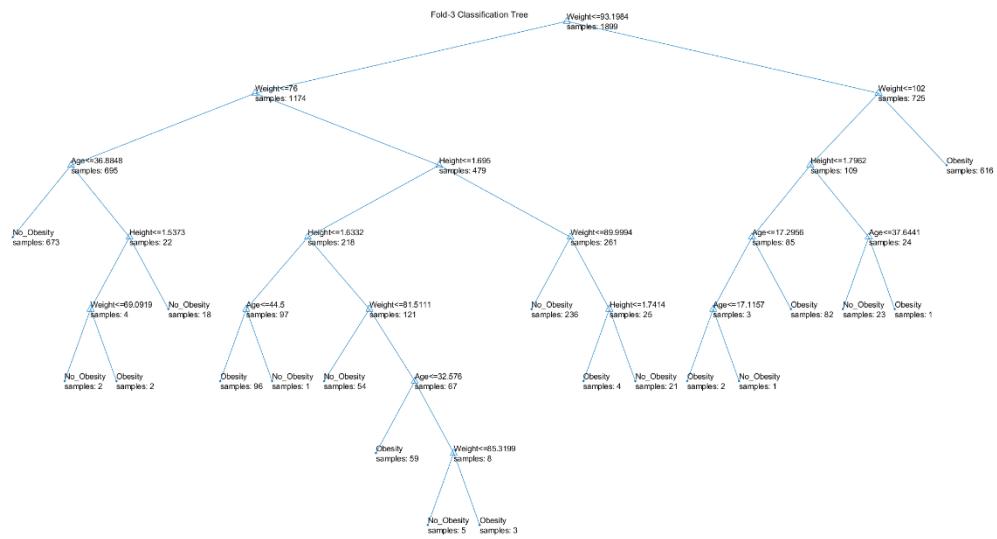
Reference

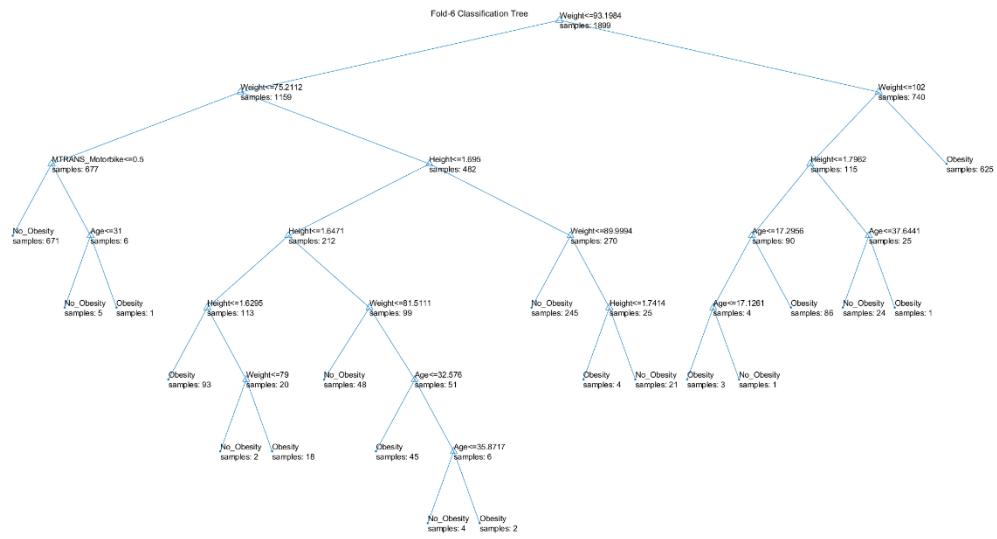
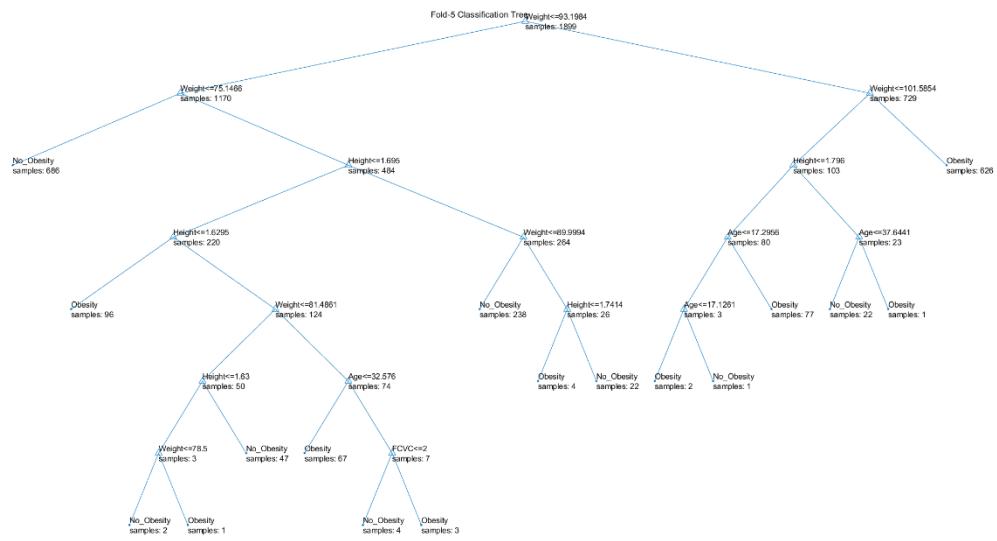
Brunk C. A. and Pazzani M J. (1991) An investigation of noise-tolerant relational concept learning algorithms. In **Proceedings of the 8th International Workshop on Machine Learning**. pp: 389-393. Evanston, Illinois. Available at: <http://www.ics.uci.edu/~pazzani/Publications/An-Investigation-MLW-91.pdf> [Accessed 25th November 2020].

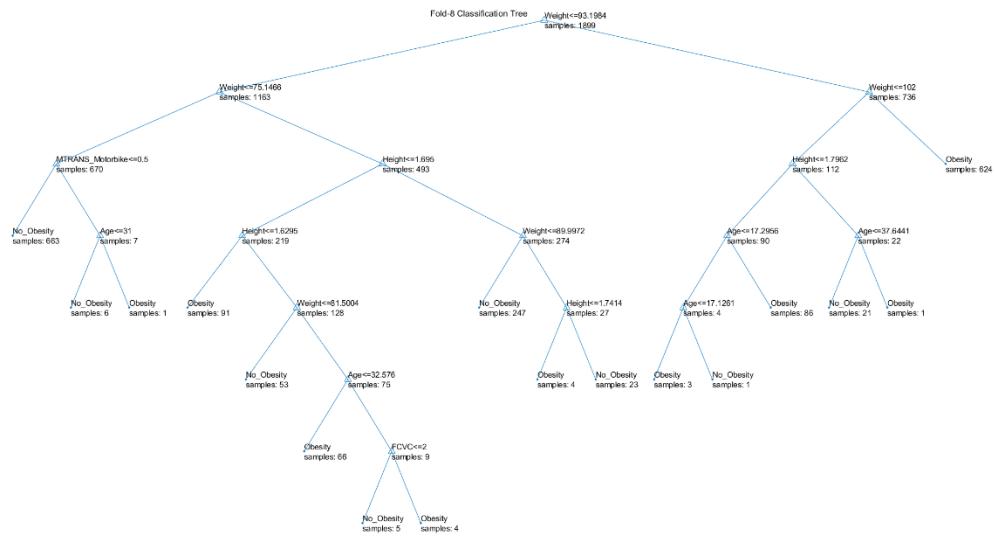
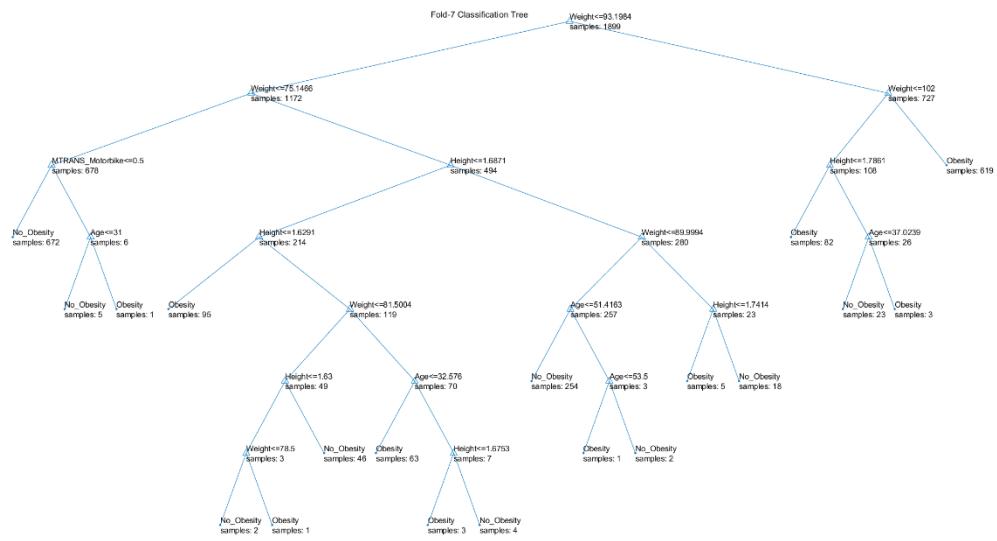
Zhang M L, Zhou Z H. (2013) A review on multi-label learning algorithms[J]. **IEEE transactions on knowledge and data engineering**. 26(8). pp: 1819-1837. Available at: <https://romisatriawahono.net/lecture/rm/survey/machine%2olearning/Zhang%20-%20Multi-Label%20Learning%20Algorithms%20-%202013.pdf> [Accessed 24th November 2020].

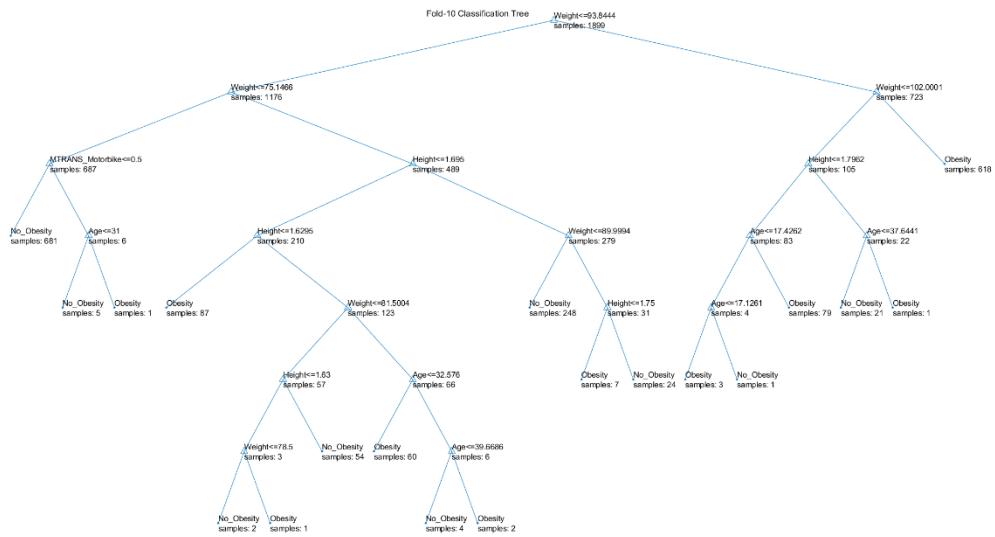
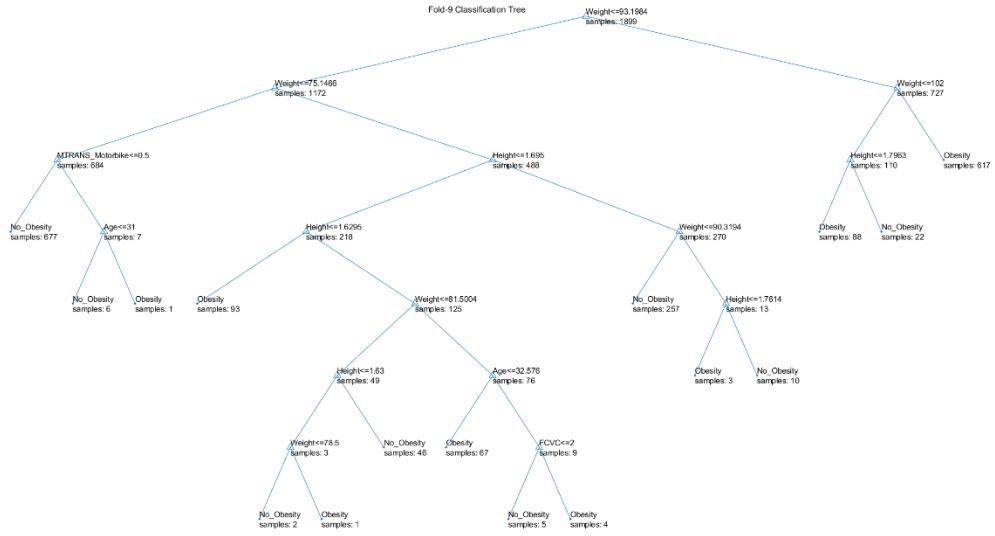
Appendix a: classification 10-fold tree drawings











Appendix b: regression 10-fold tree drawings

