

# Decision Tree Coursework Report

TszHang Wong, Quix Fan, Xuan Cai, Kaiwen Liu

October 29, 2022

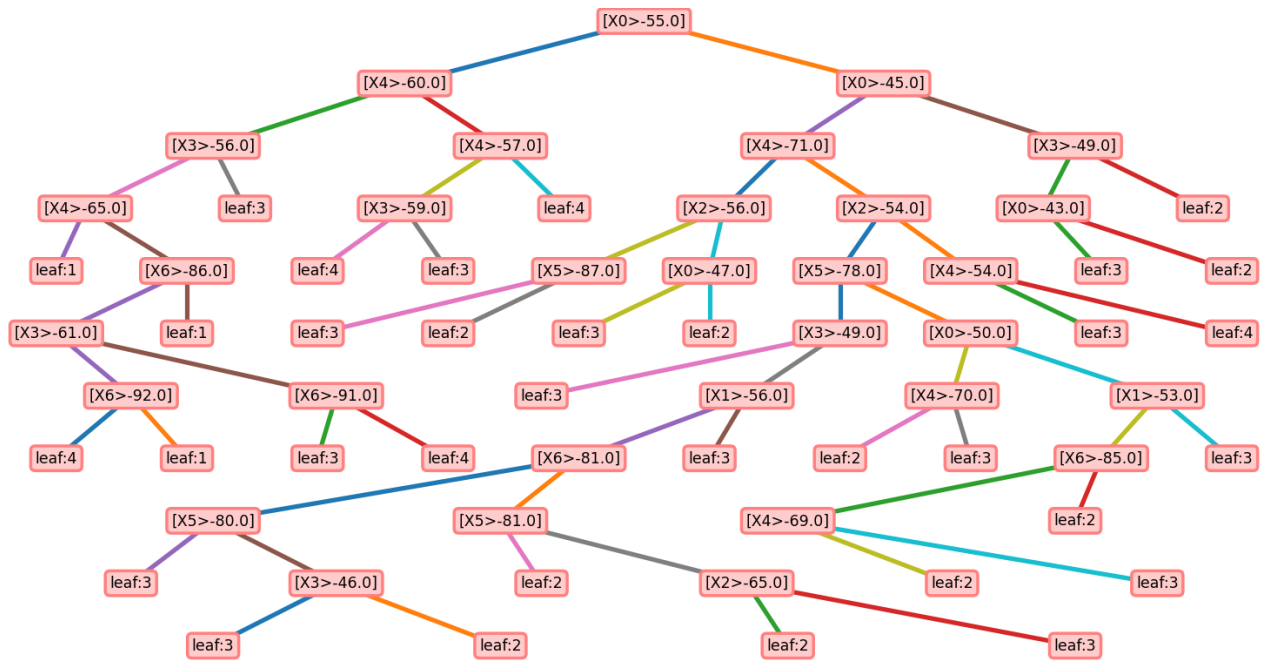


Figure 1: Tree Diagram of clean data set

# 1 Evaluation

## 1.1 10-Folds Cross validation classification metrics

	Room1 Predicted	Room2 Predicted	Room3 Predicted	Room4 Predicted
Room1 Actual	49.7	0.	0.1	0.2
Room2 Actual	0.	48.	2.	0.
Room3 Actual	0.5	1.9	47.4	0.2
Room4 Actual	0.4	0.	0.25	49.5

Table 1: Clean set confusion matrix.

	Room1 Predicted	Room2 Predicted	Room3 Predicted	Room4 Predicted
Room1 Actual	39.1	2.8	3.	4.1
Room2 Actual	3.	40.7	3.	3.
Room3 Actual	3.3	3.5	42.2	2.5
Room4 Actual	3.9	2.1	2.7	41.1

Table 2: Noisy set confusion matrix.

	Precision	Recall	F1-score	Accuracy	$\sigma_{Accuracy}$
Room1	0.9822	0.994	0.9881	0.9730	0.0103
Room2	0.9619	0.96	0.9610		
Room3	0.9556	0.948	0.9518		
Room4	0.9920	0.99	0.9910		

Table 3: Clean set metrics.

	Precision	Recall	F1-score	Accuracy	$\sigma_{Accuracy}$
Room1	0.7931	0.7980	0.7955	0.8155	0.0351
Room2	0.8289	0.8189	0.8239		
Room3	0.8291	0.8194	0.8242		
Room4	0.8107	0.8253	0.8179		

Table 4: Noisy set metrics.

## 1.2 Result analysis

Among the four rooms in the clean set, Room2 and Room3 are relatively low in precision, recall, and f1-score. These two rooms are confused as their confusing cases are the highest beside the diagonal shown in the confusion matrix. In comparison, Room4 has the highest performance among all the metrics. Regarding the noisy set, Room1 and Room4 has the least accuracy and are confused, as most false predictions lie between them in the confusion matrix. Moreover, the same performance is shown in predicting Room3 and Room2.

## 1.3 Dataset differences

In contrast to the noisy dataset, all the rooms in the clean dataset perform much better in all metrics. According to the training algorithm, a decision tree must separate the training set into subsets with a pure label to reach the optimal information gain. Hence, each noisy data is considered and captured by the model. Different labels may mix up in high-dimensional space due to noise. Thus more nodes are needed to separate them, resulting in a more complex tree, which leads to overfitting and poor generalization performance.

## 2 Pruning - evaluation

### 2.1 Cross validation classification metrics after pruning

		Room1 Predicted	Room2 Predicted	Room3 Predicted	Room4 Predicted
Room1 Actual	Before	49.3	0.	0.367	0.333
	After	<b>49.6</b>	0.	0.38	0.01
Room2 Actual	Before	0.	<b>47.9</b>	2.1	0.
	After	0.	47.62	2.38	0.
Room3 Actual	Before	0.2	1.9	<b>47.51</b>	0.39
	After	0.69	2.09	46.8	0.46
Room4 Actual	Before	0.37	0.	0.18	<b>49.5</b>
	After	0.45	0.	0.26	49.29

Table 5: Clean set confusion matrix before and after pruning.

		Room1 Predicted	Room2 Predicted	Room3 Predicted	Room4 Predicted
Room1 Actual	Before	38.3	3.43	3.22	4.03
	After	<b>44.</b>	1.3	1.4	2.26
Room2 Actual	Before	3.	39.79	4.07	2.84
	After	1.9	<b>44.04</b>	2.54	1.21
Room3 Actual	Before	2.87	3.66	41.83	3.14
	After	2.22	3.31	<b>44.08</b>	1.89
Room4 Actual	Before	3.72	2.48	3.43	40.17
	After	2.52	1.43	1.71	<b>44.13</b>

Table 6: Noisy set confusion matrix before and after pruning.

		Precision	Recall	F1-score		Accuracy	$\sigma_{Accuracy}$	Depth <sub>avg</sub>
Room1	Before	<b>0.9886</b>	0.986	<b>0.9873</b>	Before	<b>0.9708</b>	0.0144	<b>11.88</b>
	After	0.9775	<b>0.9922</b>	0.9848				
Room2	Before	<b>0.9618</b>	<b>0.958</b>	<b>0.9599</b>	After	0.9664	<b>0.0157</b>	8.29
	After	0.9580	0.9524	0.9552				
Room3	Before	<b>0.9473</b>	<b>0.9502</b>	<b>0.9487</b>	After	0.9664	<b>0.0157</b>	8.29
	After	0.939	0.9404	0.9432				
Room4	Before	0.9811	<b>0.988</b>	0.9846	After	0.9664	<b>0.0157</b>	8.29
	After	<b>0.9870</b>	0.9870	<b>0.9870</b>				

Table 7: Clean set metrics before and after pruning.

		Precision	Recall	F1-score		Accuracy	$\sigma_{Accuracy}$	Depth <sub>avg</sub>
Room1	Before	0.7998	0.7819	0.7907	Before	0.8005	0.0288	<b>18.28</b>
	After	<b>0.8688</b>	<b>0.8980</b>	<b>0.8831</b>				
Room2	Before	0.8062	0.8006	0.8034	After	<b>0.8813</b>	0.0288	13.42
	After	<b>0.8793</b>	<b>0.8862</b>	<b>0.8828</b>				
Room3	Before	0.7960	0.8123	0.8041	After	<b>0.8813</b>	0.0288	13.42
	After	<b>0.8855</b>	<b>0.8559</b>	<b>0.8704</b>				
Room4	Before	0.8003	0.8066	0.8034	After	<b>0.8813</b>	0.0288	13.42
	After	<b>0.8918</b>	<b>0.8862</b>	<b>0.8870</b>				

Table 8: Noisy set metrics before and after pruning.

## 2.2 Result analysis after pruning

A slight decrease in average accuracy is shown after pruning for the clean set because of the loss of non-noisy information captured within the train set. In contrast, the accuracy increased considerably in the noisy dataset after pruning. Pruning eliminated the captured noise by replacing nodes with the dominant label set when accuracy remains reasonable, resulting in a less overfitted decision tree with a higher bias and lower variance than before.

## 2.3 Depth analysis

The tree trained on the noisy dataset is shown to have a much larger depth than the clean dataset because more nodes are needed to separate data under the noisy data pattern. Pruning resulted in a tree with less average depth on both datasets, with prediction accuracy increasing in the noisy dataset and a mighty decrease in the clean set. Overall, the tree's generalization performance increases with less complexity and smaller maximum depth.