# Data Partitioning and Modeling

The data was partitioned into train and test datasets.
The train data set was used to create the decision tree model.
The trained model was then applied to the test dataset.
This is important because we want to make sure our classifier is not overfitting or underfitting.

When partitioning the data using sampling, it is important to set the random seed because we want to obtain the same training and test datasets to train and test the decision tree algorithm throughout the participants in the course.

The data was partitioned into 60% train and 40% test datasets. The training data is used to build the classification algorithm, whereas the test data is used to evaluate the classifier on new, unseen data to make sure our classification is not overfitting or underfitting. For this project, we selected the stratified sampling method which divides the datasets into separate groups referred to as strata and a probability sample is drawn from each group. When partitioning the data, it is important to set the random seed to a certain number which in this case is 1466016757670 so that we obtain the same training and test datasets to train and test the decision tree algorithm.

A screenshot of the resulting decision tree can be seen below: