# Data Partitioning and Modeling
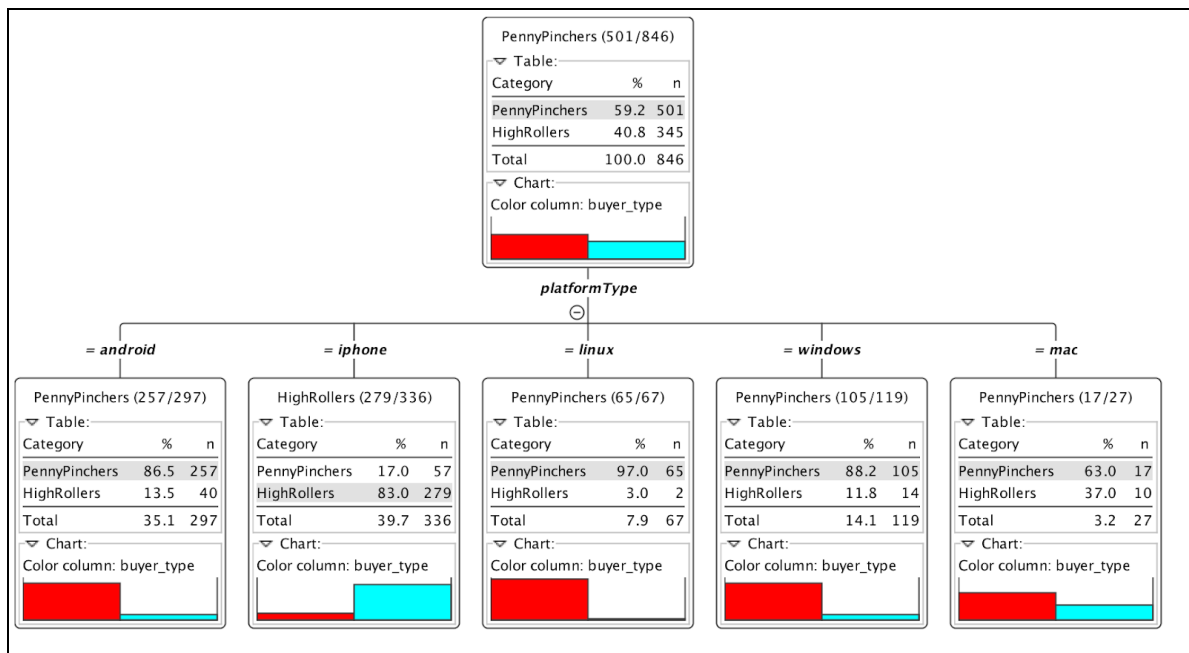
The data was partitioned into train and test datasets.
The 60% training data set was used to create the decision tree model.
The decision tree model was then applied to the remaining 40% test dataset.
To split the data into a training and test datasets is important because it allows us to build a model based on the training data and then apply this model to the remaining rows (test data set) to measure the accuracy of classification of our model.

When partitioning the data using sampling, it is important to set the same random seed because we can have consistent (same) data partitions when partition node is executed.

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





[root]:  class 'PennyPinchers' (501 of 846)

[platformType = android]:  class 'PennyPinchers' (257 of 297)

[platformType = iphone]:  class 'HighRollers' (279 of 336)

[platformType = linux]:  class 'PennyPinchers' (65 of 67)

[platformType = windows]:  class 'PennyPinchers' (105 of 119)

[platformType = mac]:  class 'PennyPinchers' (17 of 27)