Homework 2 Analysis

Based on the decision tree analysis for predicting Parkinson's disease, the tree depth reached five nodes. In pursuit of higher accuracy, we tried changing the max_depth function and we found out the tree will never go over 6 nodes even if the max depth would be set over 6. In contrast, when the depth was set to lower than 5, the training accuracy decreased.

The tree created 13 leaf nodes indicating the tree's decision outcomes. Only 7 out of 22 attributes were used, highlighting their higher relevance in predicting the disease. Most likely, the other attributes had a lower information gain for their current data, therefore not being that relevant.

 The confusion matrix revealed a tendency towards false positives over false negatives, which, while not ideal, is preferable in medical diagnostics to avoid missing potential Parkinson's cases. Our tree is more likely to misclassify a healthy person as a sick person rather than the other way around. It is preferred because our main focus is predicting the sick cases rather than preventing fake positive diagnoses. A sick person should start treatment as soon as possible, especially considering the progressive nature of Parkinsons.

The substantial accuracy gap between training (100%) and testing data (60%) would usually suggest overfitting, but in our case we can't assume that, taking into consideration our really small testing sample. If the gap consists with bigger testing samples, overfitting would be a fair assumption.

Improving classifier performance could involve refining the tree's complexity, possibly by adjusting the max depth or the minimum number of samples required for a split, and expanding the testing dataset for a more robust evaluation. Additionally, pruning the tree would also help remove branches that do not help the efficiency. Changing the random state of the tree

would result in more constant accuracies and a cross validation would create a better validation

system before using the testing data.