**Report**

**Feature Detection/Methodology & Findings:**

I took the provided data and the data labelled earlier and listed features in a single csv. Added features and attributes to the same file and loaded the file in weka for visualization with the feature set. I tried to take the attributes with truth/False values as I estimated that the classes are distributed in two. I was looking for the feature(s)/set of feature(s) that gives me the similar distribution of data as the number of classes.

Extracted Feature List: [length, vovel_count, consonent_count, start_with_vovel, Len>4, even_len, first_vovel_position, first_consonent_position, ends_with_vovel, second_is_vovel]

I checked repeatedly with many different combinations. But the distribution was nowhere near the number of classes for combinations. I also tried different classifiers with different feature combinations but most of them gave 50-100% error.

I found out on manual analysis that the data I included was either differently mapped or wrong. So I took the original dataset provided in the assignment google drive and grouped the class instances together to check the similarity and I easily found out that the feature seems to be the second vowel alphabet.

Then I checked the features separately. Adding closest feature distribution map here:
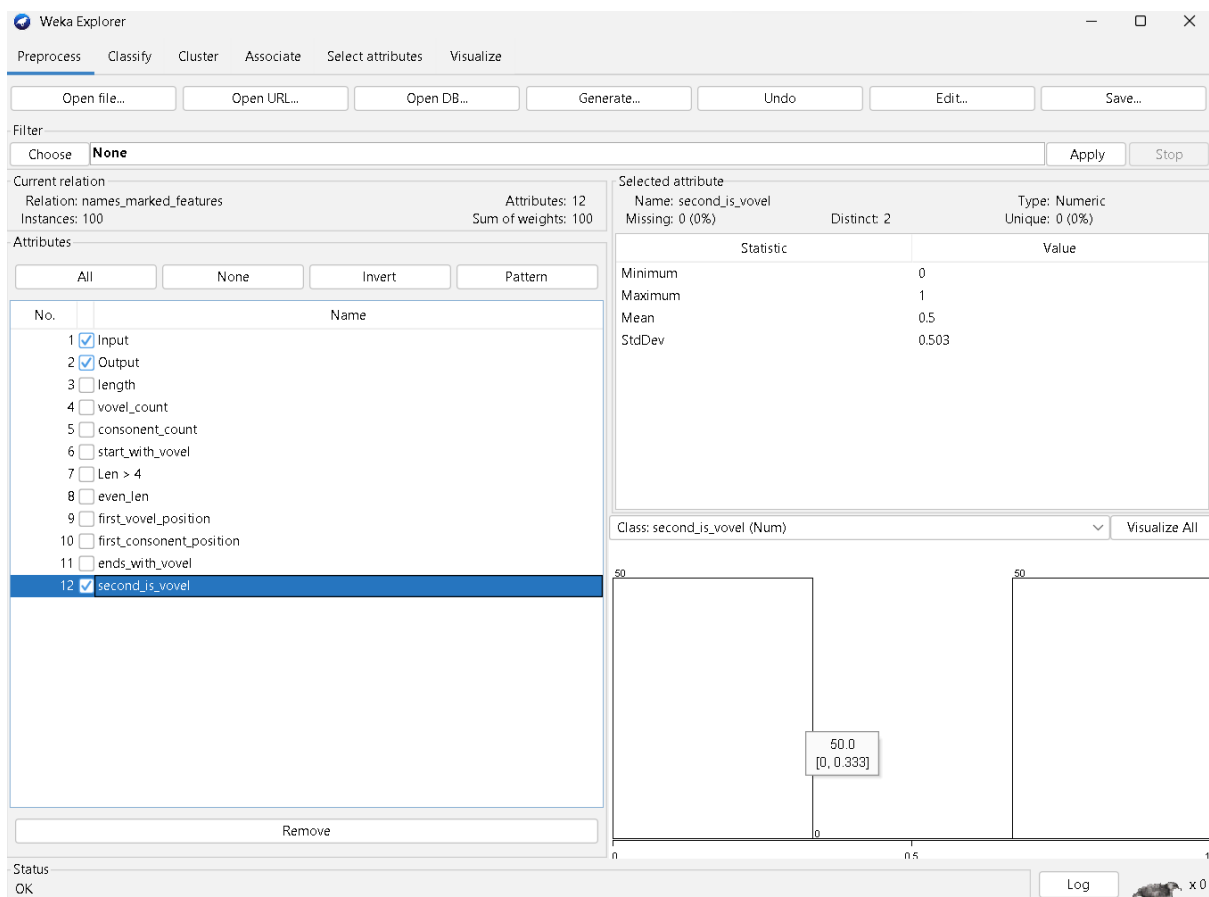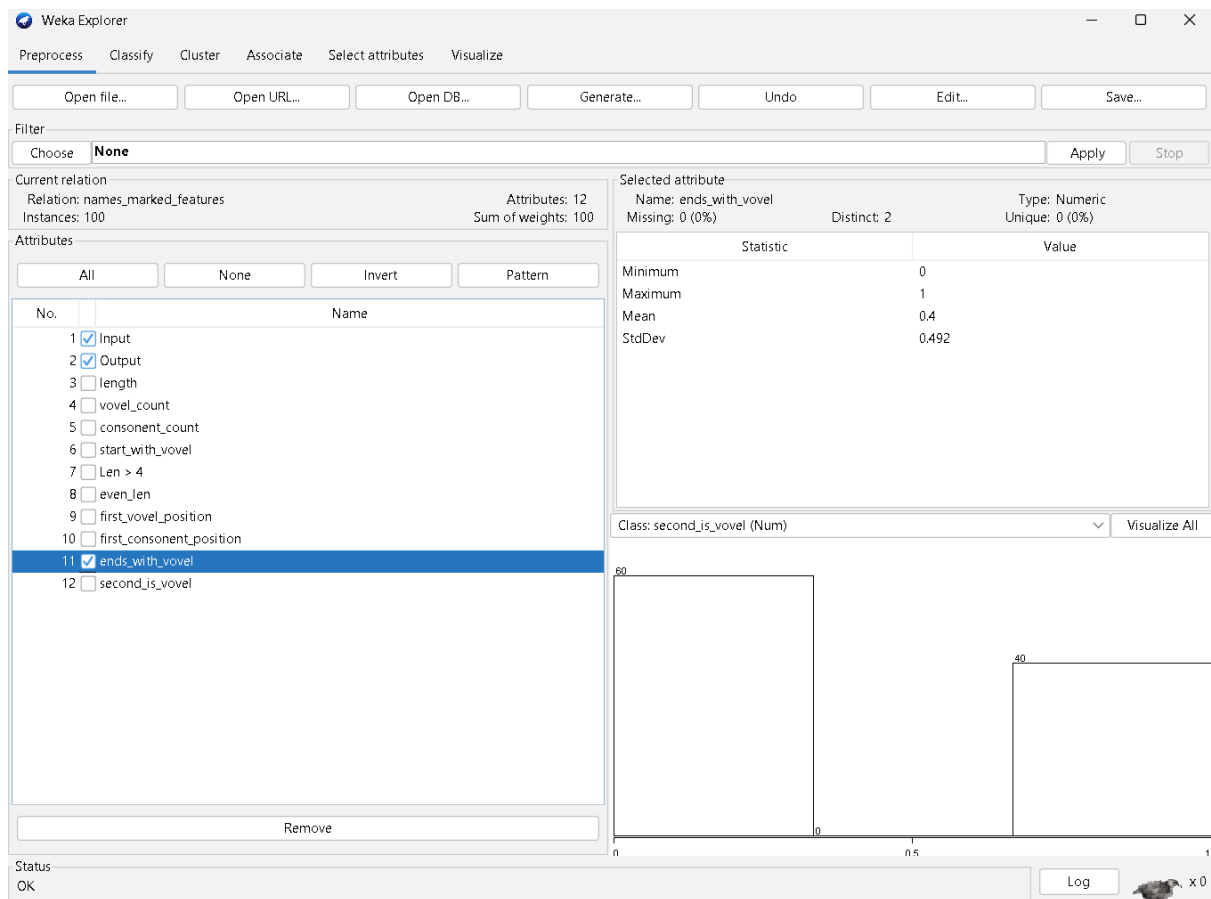


*Figure 1 The correct Feature as per distribution*

*Figure 2 Checking for features*

For my analysis, I checked and analyzed other feature data as per my own understanding and shortlisted some other features to experiment with. Although the results were expected, but there were some other relations as well, like first_vowel and consonant_count had distribution near that of results.

**Results:**

Since the feature detected was right, the classification result should be 100%. And with general decision tree classifiers, most of them gave 100% result.

Attaching the result of final tries of the classification algorithm implementation below:



*Figure 3Final Result from J48 algorithm*

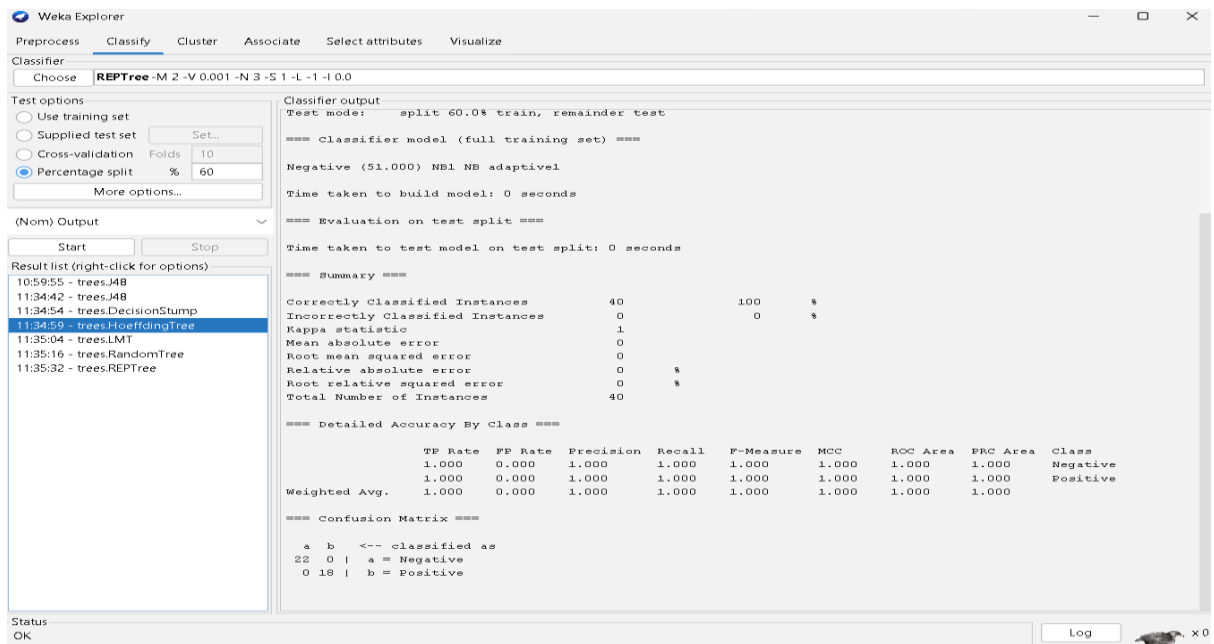Other algorithm with 100% results are as below:



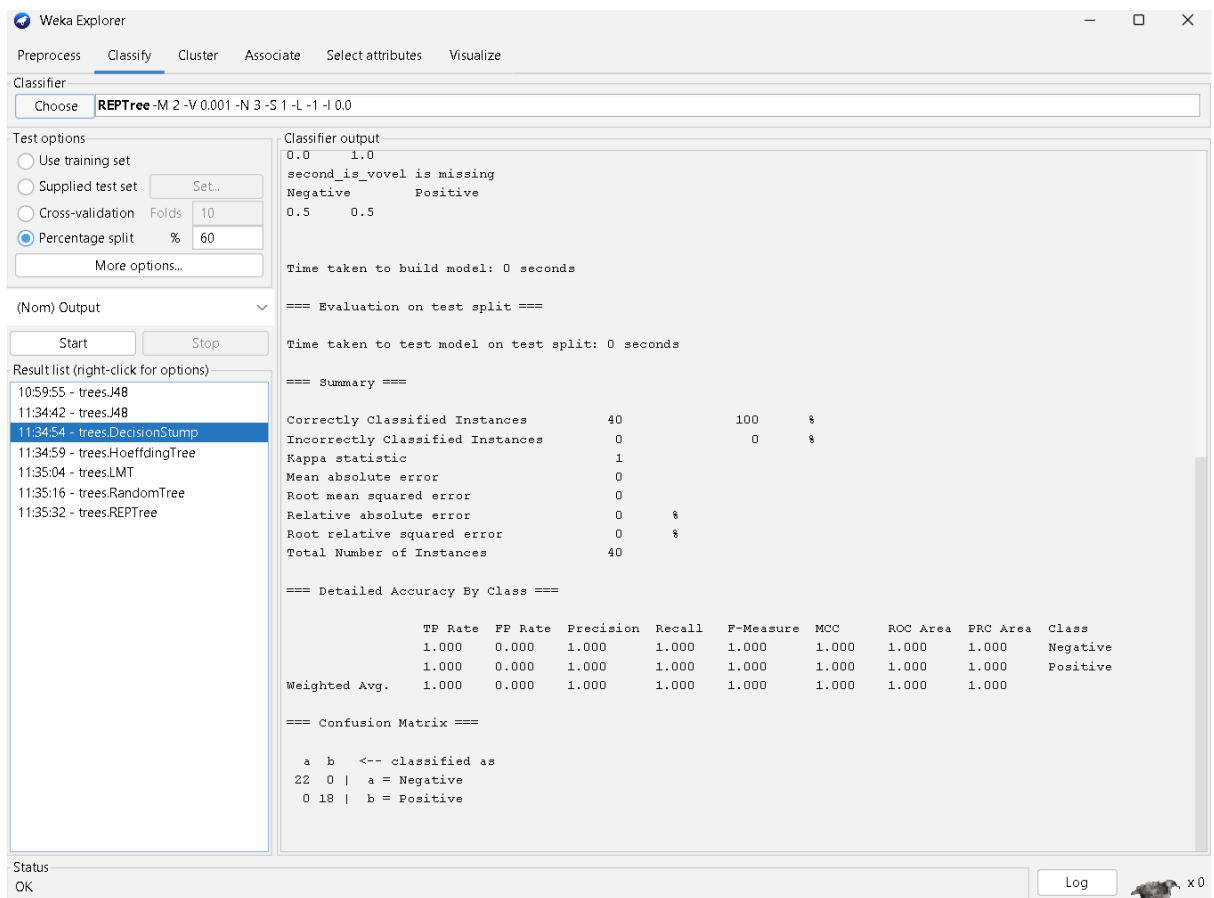*Figure 4Hoeffding Tree Results*



*Figure 5 DecisionStump results*

Other algorithms like LMT, REPTree and RandomTree gave 26-100% errors in different tries and split