



Student ID Number: 19075153

Assignment 1 Part B

Semester 2 2021

Student Name: Bernard O'Leary
Student ID: 19075153

PAPER NAME: Data Mining and Machine Learning

PAPER CODE: COMP809

Due Date: Friday 10 Sep 2021 (midnight)

TOTAL MARKS: 100

INSTRUCTIONS:

1. The following actions may be deemed to constitute a breach of the General Academic Regulations Part 7: Academic Discipline,
 - Communicating with or collaborating with another person regarding the Assignment
 - Copying from any other student work for your Assignment
 - Copying from any third-party websites unless it is an open book Assignment
 - Uses any other unfair means
2. Please email DCT.EXAM@AUT.AC.NZ if you have any technical issues with your Assessment/Assignment/Test submission on Blackboard **immediately**
3. Attach your code for all the datasets in the appendix section.

Study Area I (Dataset is bank.csv use the Bank.zip)

(a) Pre-processing

All of the libraries necessary for the assignment were imported in the first step, including those necessary for pre-processing the data.

Preprocessing data required the application of semi-colon delimitation, removal of any records that were classified "unknown" and turn all parameters into categorical (numerical) values so that they can be processed by the various classifiers that we are using for this assignment.

(b) Top five most influential features

The top five most influential features were selected by creating a method that will test the accuracy of models created using a range of different feature extraction processes. The processes used are:

- Extra Trees Classifier
- RFE (Recursive Feature Extraction) using Logistic Regression
- Chi-squared test
- Principal Component Analysis (PCA)

The features were then ordered in order of influence, or captured using the algorithm's built-in feature selection capacity. The "transform" method on the fitted model object is then used to enable the top five features to be selected from the result.

The best performing approach was RFE. The RFE process has in-built ability to refine the feature-set to five most influential features. The process for feature selection using RFE is as follows:

```
# Feature Extraction with RFE
model = LogisticRegression()
rfe = RFE(model, 5)
fit = rfe.fit(X, Y)
features = fit.transform(X)
pred_features = features[:, 0:5]
```

Once the top five features have been identified, they are run through a "get_accuracy" method which creates a MLPClassifier model and returns an accuracy score. These scores are then compared. As above, the best scoring approach was RFE.

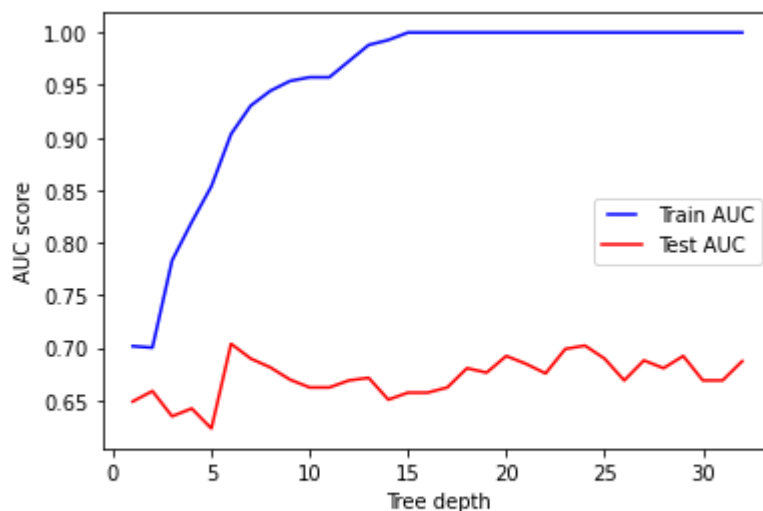
(c) Decision Tree algorithm

The two parameters that proved to be most influential in ability to reduce the number of nodes in resulting tree were the min_samples_leaf and min_samples_split parameters. The best performance was obtained when min_samples_split=0.3 and min_samples_leaf=0.2.

The min_samples_leaf parameter produced smaller trees with an optimised value, however neither parameter seemed to have any effect on the accuracy of the model training process, despite reducing the number of nodes in the tree.

The following images show the train and test AUC performance for each different type of five parameters tested. The resulting model accuracy tested with 10x cross-validation and number of nodes in the resulting tree is listed. An optimised model was trained for each parameter to test for accuracy based on what appeared to be an optimal value for the variable as is shown by the charts produced (i.e. values that appear to reduce the risk of under/over-fitting).

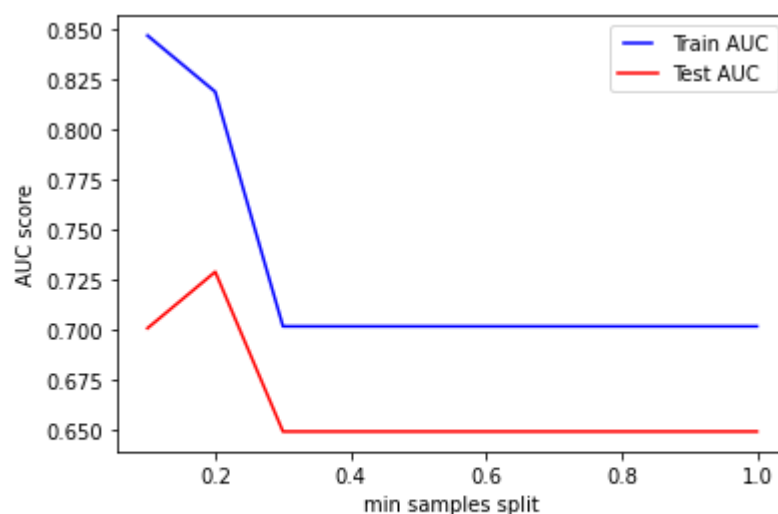
max_depth=3



0.8301572897761644

15

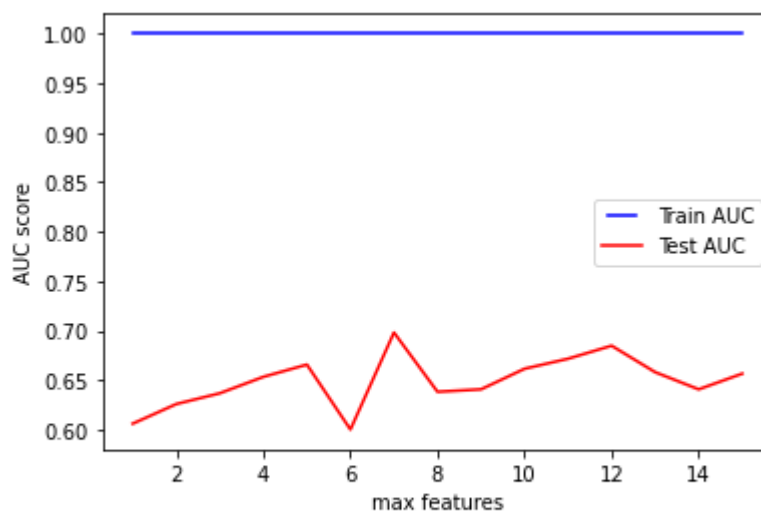
min_samples_split=0.3



0.8196612220205687

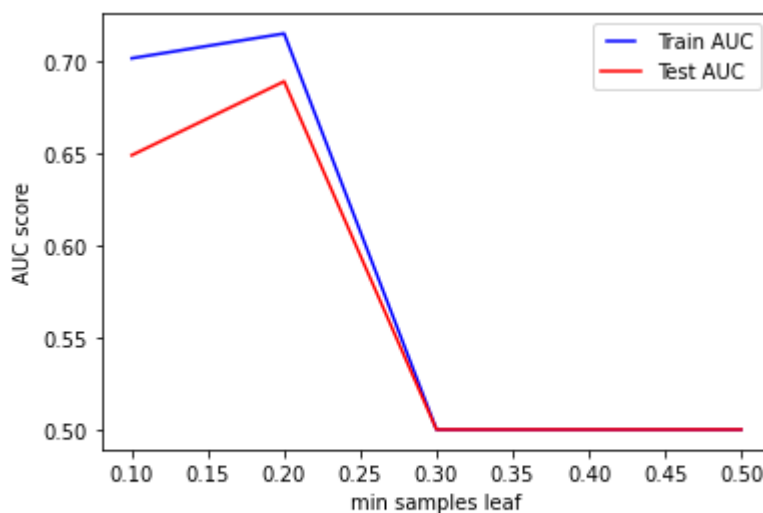
13

max_features=16



0.7968844525105868
151

min_samples_leaf=0.2



0.7827888687235329
7

(d) Describe the role of the two parameters

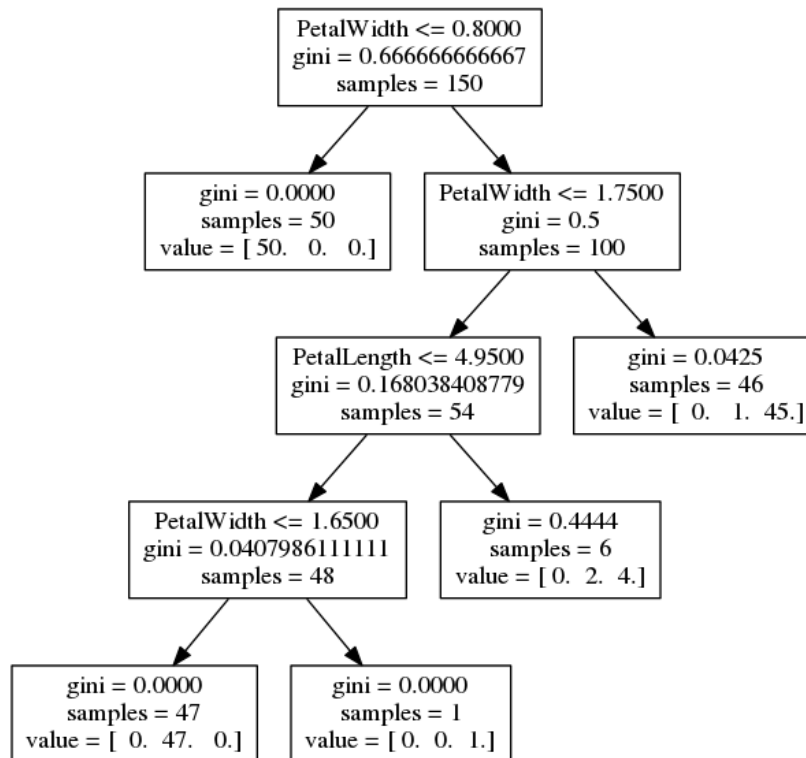
As can be seen from the above results, the best performing parameters and values for the parameters as obtained when `min_samples_split=0.3` and `min_samples_leaf=0.2`.

`min_samples_split` defines the minimum number of samples that are required to be analysed before the model generation algorithm will elect to split an *internal* node of the tree. For example, if `min_samples_split=2`, then a minimum of two observations need to be considered to split the node.

`min_samples_leaf` parameter works the same way, but for external (leaf) nodes. The difference between this parameter and `min_samples_split` is `min_samples_leaf` specifies a

minimum number of samples in a leaf, while `min_samples_split` can create as many sub-leaves as it's setting allows, dependent on the value specified for `min_samples_leaf`.

The following diagram illustrates the behaviour and subtle differences between the two parameters:



This tree was constructed against the Iris dataset with `min_samples_split=10`. We can see that the internal nodes of the tree have no samples less than 10, whereas there are two leaves that considered less than 10 samples. `min_samples_leaf=1` for this tree.

Documentation for the parameters is included in the Decision Tree Classifier documentation – here: <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

Depending on the size, shape and nature of the content of the dataset, the parameters that had an effect here may not have as strong of an effect as they have had here. For example, for large datasets such as for image processing or audio processing applications, performance of these parameters might be outperformed in terms of influence on the model by adjustment of some other parameters.

(e) Examine the Confusion Matrix

The final score, number of nodes and confusion matrix with and without normalisation are provided below. We can see that of the observations considered, we have 87% true-positives and 51% true negatives. False positive is sitting at 49% whereas false negatives is sitting at 13%.

```
0.7827888687235329
```

```
7
```

```
Confusion matrix, without normalization
```

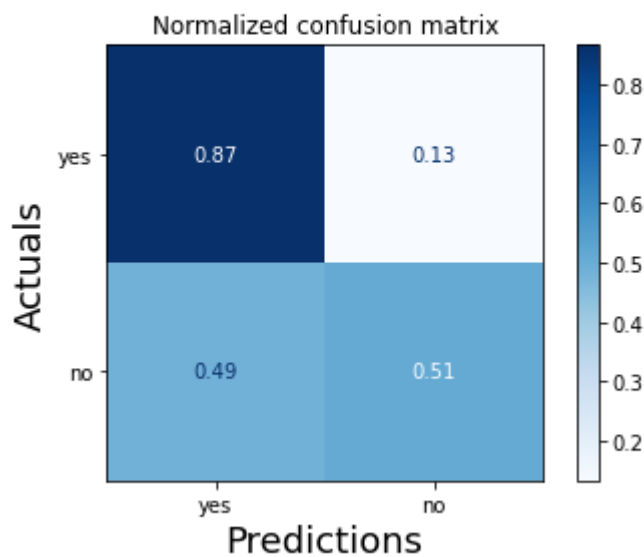
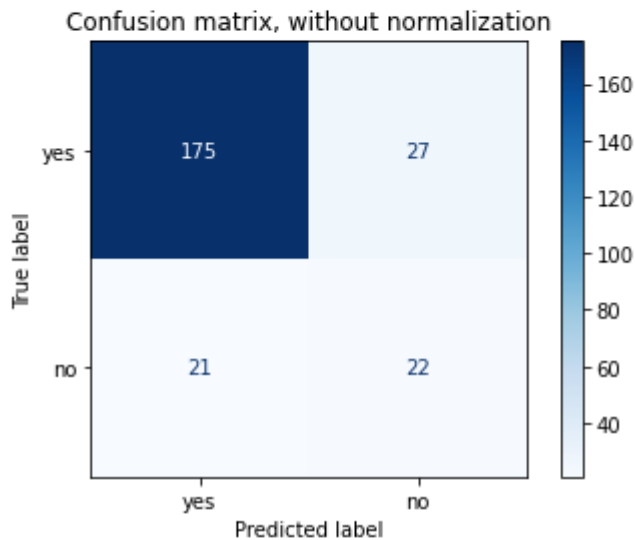
```
[[175  27]
```

```
 [ 21  22]]
```

```
Normalized confusion matrix
```

```
[[0.866 0.134]
```

```
 [0.488 0.512]]
```



The model was trained and tested with 10x cross-validation. The overall accuracy of the model is reported as 78%, which is significantly less than the percentage of true positives, but is influenced by the low number of true negatives. Our equation for accuracy is as follows:

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

It is likely that the number of “yes” results for outcomes is larger as it is in the interest of banks to have customers make term deposits, so the marketing campaign conducted might have been very effective. Also this dataset might have been taken from a time when

interest rates for term deposits were unusually high, making it a comparatively good investment. Due to globally low interest rates for lending, this is not necessarily the case nowadays.

Study Area II (Dataset is Autism-Child-Data.arff)

(a) Describe the autism spectrum disorder

The Autism-Spectrum Disorder (ASD) refers to and quantifies autistic traits in adults [1]. Because Autism is a condition that influences aspects of behaviour, such as high neuroticism, low extraversion and low agreeableness [2], it can lead to difficulty succeeding in environments such as educational institutions that do not cater for variances in behaviour that may be exhibited in someone who is identified as being higher than average on the Autism Disorder Spectrum. On the other hand, there is evidence to suggest that individuals who are identified as being higher than average on the Autism Disorder Spectrum, may favour vocations such as engineering, physics or mathematics [2], which allows education institutions to support children who score higher than average on the Autism Disorder Spectrum by providing options to focus more in these areas. It is therefore of interest to be able to identify individuals who may be on the Autism Disorder Spectrum earlier rather than later.

The test for ASD is called the Autism Spectrum Quotient (AQ). The intention of the AQ is to provide a quick and easy way to quantify how many autistic traits an adult has [1]. The AQ is a self-reporting test that consists of a 50-item questionnaire [2]. The questionnaire is split evenly into five different types of questions covering social skill, attention switching, attention to detail, communication and imagination. Answers to questions are designed so as an answer can be classified as either autistic or non-autistic. Questions that may be included are for example “I find it hard to make new friends”, or “I am fascinated by numbers”. Although the AQ test is designed originally for adults, there is no significant difference in its ability to indicate where an adolescent is on the Autism Disorder Spectrum [1].

[1] Baron-Cohen, S., Hoekstra, R.A., Knickmeyer, R. et al. The Autism-Spectrum Quotient (AQ)—Adolescent Version. *J Autism Dev Disord* 36, 343 (2006).

[2] Elizabeth J. Austin, Personality correlates of the broader autism phenotype as assessed by the Autism Spectrum Quotient (AQ), *Personality and Individual Differences*, Volume 38, Issue 2, 2005, Pages 451-460

(b) Identify the top five significant features

The same approach for feature selection was applied in this section as it was for the first section. The top five features were identified using the Extra Trees Classifier method. This method was selected because of the four different feature selection methods that were tested, Extra Trees Classifier performed slightly better than the rest. A typical result was as follows:

Accuracy score of our model without feature selection : 0.05

Accuracy score of our model with chi square feature selection : 0.15

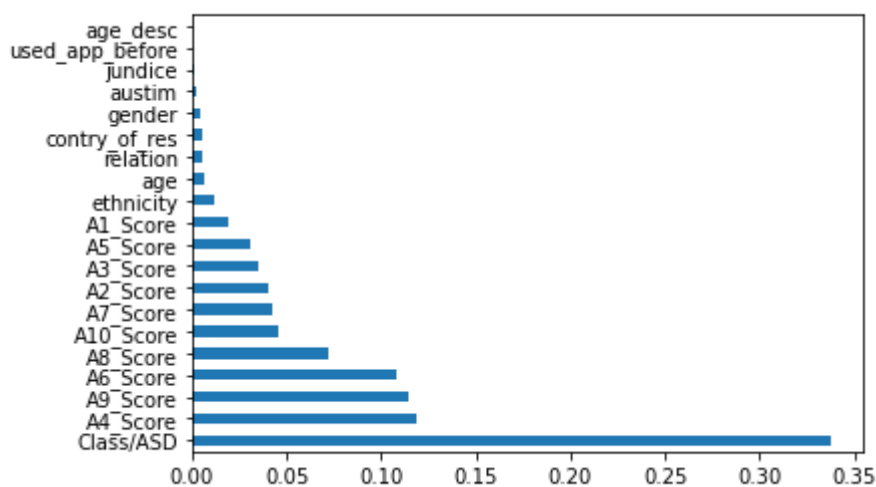
Accuracy score of our model with RFE selection : 0.15

Accuracy score of our model with PCA selection : 0.15

Accuracy score of our model with Extra Trees selection : 0.16

Documentation for the method can be found here: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesClassifier.html>

Several data cleansing steps needed to be applied prior to application of the Extra Trees Classifier. Cleansing included changing the data to be categorical and changing the imported to remove quite marks around all of the data that were imported so that it could be transformed to categorical. Finally rows that have negative values were removed from the dataset. The result of applying the Extra Trees Classifier is illustrated in the following chart. In this case, the five highest scoring features were A8_Score, A6_Score, A9_Score, A4_Score and Class/ASD. The results varied each time the process was run, but Class/ASD always shows up as most influential, with a combination of other features thereafter.



The independence assumption states that each feature is conditionally independent of any other feature. This is a requirement of the mathematical theory behind the Naïve Bayes algorithm, however in practice, it will often not always be the case that all features are entirely statistically independent of each other. If a large numbers of features are dependent on each other though, accuracy may drop.

Because there are only 10 questions asked by the AQ test given in the study, and we are not give the details of what the questions were (just the index) it is difficult to just by looking at the questions how likely it is that they are properly independent of each other. It is likely that the final score (Class/ASD) is not fully independent from any of the individual the questions though.

(c) Naïve Bayes algorithm

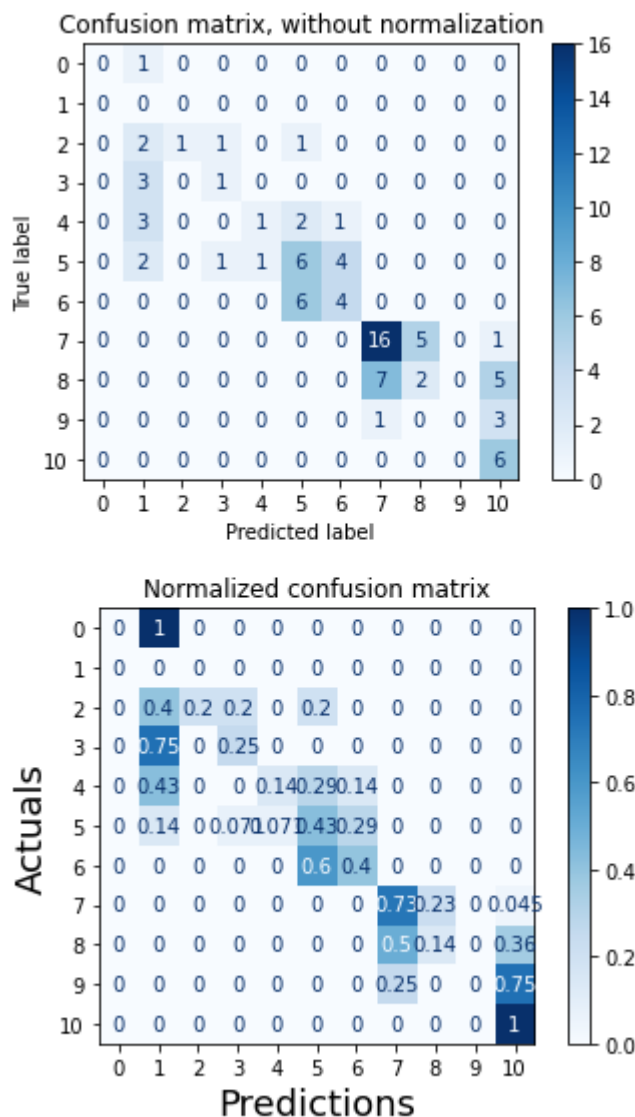
The accuracy for the NB algorithm was consistently sitting at about 35% accuracy. This indicates that the model would be able to product the correct AQ score (a number between 1 and 10) about 35% of the time. This is not a great result, and the model is therefore not considered a good predictor.

This might be because of the smaller number of features that are included in this version of the AQ model, which is for the intention of using with children, to make it easier for them to complete. The original AQ model has 50 features in it, whereas this modified one for use with children has only 10. Where the independence assumption is not completely

met by all the features for example, with less feature, the issues caused by this might be more prominent.

A way to reduce this effect might be to include more data, rather than the 292 records included for this study. With more data, the accuracy might improve.

We can see from the confusion matrix that there is the start of some strong correlation between features through the horizontal of the diagram, which is what we would want to see to indicate a high proportion of true positives and true negatives, but for fault positives and false negatives, it looks like there is just not enough data for the model to give a reliable estimate for how often these situations might occur.



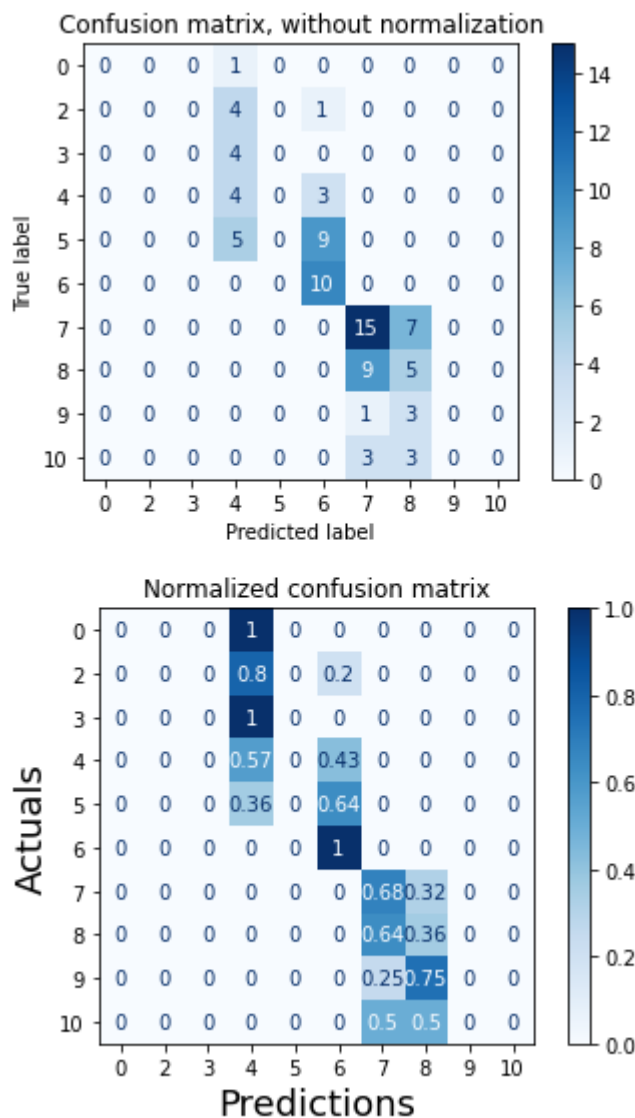
0.3580952380952381

(d) Decision Tree Classifier algorithm

Result of the Decision Tree Classifier algorithm was significantly less accurate than the NB, although clearly neither is a good predictor, the DTC algorithm did worse than the NB one. Accuracy was consistently reported as being about 25% - about 10% less accurate than for NB. This indicates that the model would be able to product the correct

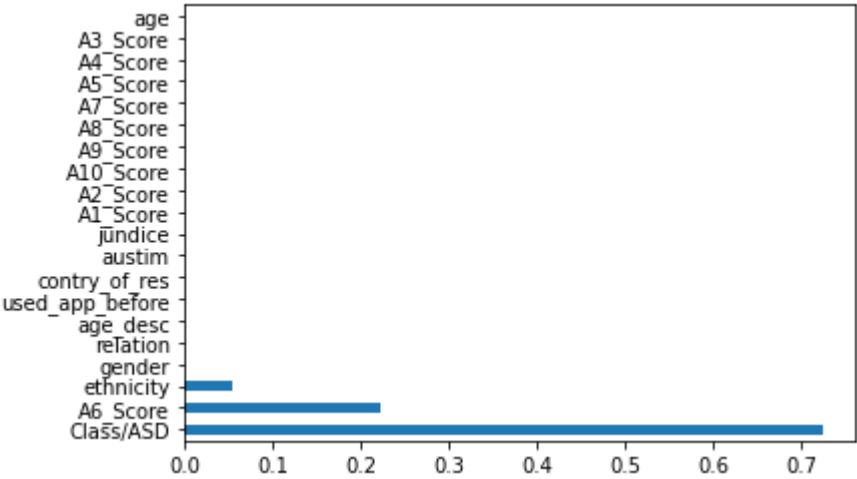
AQ score (a number between 1 and 10) about 25% of the time. This is not a great result, and the model is therefore not considered a good predictor.

As we can see from the confusion matrix, there is some clustering to the mid-diagonal of the chart, but the data are clumped and less well distributed, indicating less predictive capability than for NB just by looking at the quality of the confusion matrix.



0.27809523809523806

The top five features for the DTC model were Class/ASD, A6 _Score, ethnicity, gender, relation. The influence of each feature is less well distributed, however Class/ASD in this case is also still the feature with the most influence, same as for NB. The DTC seems to place less significance on the binary features in the dataset, i.e. the questionnaire answers (y/n) and places more emphasis on the categorical features such as ethnicity. This is possibly because the NB model is better suited to binary data.



Appendix