

Learning from data Coursework

December 2021

Contents

1	Introduction	3
2	Data Set	3
3	Training Summary	5
4	Key Findings and Insights	8
5	Next Steps	10
6	Appendix 1	13
7	Appendix 2	19

1 Introduction

The main objective of this work is to produce models that can predict the quality of a wine based on its chemical composition. This prediction will take the form of the models assigning a wine with a number for 0 to 10 as a ranking, therefore making this a classification problem then is further analysed to provide predictive insight. As the problem is classification in nature the two models chosen are multilayer perceptron (MLP) and convolution neural networks (CNN).

To gain a full understanding of the data and its predictions the analysis will be split into multiple subtasks: An outline of the data is required, what are the possible relations between different chemical components? Furthermore is the relation between each individual chemical component and the related score, this is to identify if any particular component is more influential than the others and will help in the generation of chemically superior wines that can then be assessed with the models. The data as a whole must be analysed in comparison with the ranking as this may produce further interactions that are not seen individually between components. These steps should provide insight into the components of the wine in order to predict what levels of each are required to create superior wines and produce models that are able to predict a wine's quality.

There is one significant potential issue with the data and that is the nature of the quality statistic. Due to this having been defined subjectively by one or more people it. As different people have different tastes, the data may not follow a perfectly logical pattern and may not produce the best results. Therefore the models may not produce accurate predictions.

2 Data Set

The data set [1] comprises of a [12, 1600] arrangement of data that are all numerical, apart from the titles of course. There are 11 features to the data set, these are; fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates and alcohol. There is also one outcome variable that is the quality. All values are positive numbers without units as for this purpose the units are not required, and the quality is between 0 and 10. However when it comes to quality the data is not balanced as there are only 6 of the 10 possible classifications present and they are all in the middle. See Figure 1.

This dataset is one of the few cases where there is complete data, without any missing values and complete, unified data. Therefore very little feature cleaning or engineering is required. However in the future as more data may be collected complete data coverage cannot

be guaranteed and therefore must be guarded against. If the data was just set to 0 it would drastically effect the models so another option is needed. If the data is received with a quality value then all of the wines in the same quality category can have their features averaged and the average is then used to fill the missing data. This will slow down the pre-processing however it will provide a much more reliable standard for what the value should be improving the accuracy of the model.

Further to this, analysis of how the relationship between the features and the outcome variable shows there are some features that effect the outcome variable more than others. For each feature I have normalised its data to produce a standardised gradient, allowing me to accurately analyse the data produced, the gradients are all to three significant figures. See appendix 2 for all of the graphs produce. The first feature is fixed acidity which produces a noticeable, yet small positive gradient of 0.017 showing that it does have an impact on the result however it is not huge. Volatile acidity and citric acid have similar, yet opposite effects on the quality of the wine with gradients of -0.055 and 0.055 respectively. These values are relatively large compared to the other features showing they are areas of interest in predicting the quality of a wine. However some areas appear to not hold much value in predicting the quality, this is the case for residual sugar with a very small positive correlation of 0.002 showing it has some effect but very small. The chlorides have a similar, yet negative effect to the fixed acidity having a gradient of -0.012, therefore these should be reduced when creating a quality wine. There is an interesting difference between two closely related features. The levels of free sulfur dioxide seem to produce a small negative effect with a gradient of -0.009 and very little effect on the quality, however the total amount of sulfur dioxide produces a large effect on the quality with a gradient of -0.026. With both of these values being negative and free sulphur being derived from the total sulphur, a simple area to work on to improve the quality is to reduce the total sulphur levels in the wine. All of the features have some effect on the quality except for the density that produces a gradient of 0.0 with three significant figures and therefore does not seem to matter when creating wine. The pH is another feature with very little effect on the quality of the wine, with a negative gradient of -0.003 and therefore is not a vital feature to check. The final two features are much more important to the quality of the wine, with the amount of sulphates having a positive 0.026 gradient and the alcohol level having a positive 0.042 and therefore are vital when creating a quality wine.

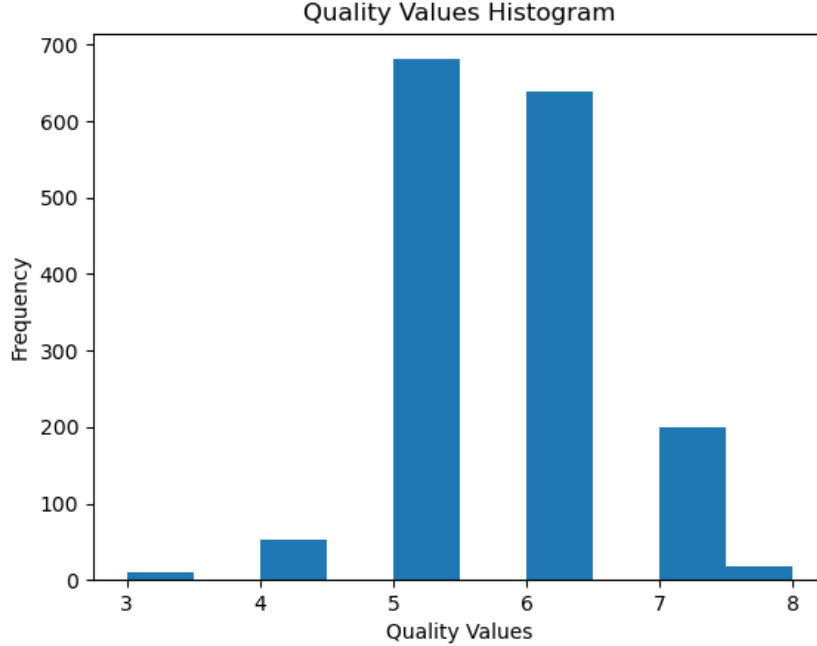


Figure 1: A histogram showing the distribution of quality data in the data. This shows a large tendency towards the middle numbers and very few highs and lows.

3 Training Summary

As stated previously the two models being used to analysis the data are Multilayer Perceptron models and Convolution Neural Networks, all of these tests are run with the full values of the wine at once as the model may provide some insite or find some connections that doing it individually would not find. The implementation of the MLP provides many different variables to change in an attempt to provide the best results, these are the size of the hidden layers, maximum iterations, the activation function, the solver and the learning rate. The initial implementation has 5 hidden layers, 100 max iterations, the relu activation function, a Stochastic Gradient Descent solver and inverse scaling learning rate. When run 50 times with a test/train split of 0.2 it produces an average mean accuracy of 0.34 which is quite low and therefore requires a change in variables. Increasing the size of the hidden layer to 10 and the max iterations to 1000 improves the mean accuracy to 0.41 with yet another increase of the hidden layer size to 12 providing no additional benefit with

a mean accuracy of 0.412 again. therefore the previous noticeable increase may be due to the max iterations, however increasing this to 1500 provides very little improvement, with a mean average of 0.43, this appears to be because the training loss does not improve and the model stops its self after 10 consecutive epochs if it does not improve by more than 0.0001. This leads me to believe the learning rate needs to be changed in order to improve the model, so changing the learning rate to constant provides a drastic benefit, raising the average mean to 0.51 with a adaptive learning rate improving on this very slightly at 0.52. The activation function is a crucial part of any model and therefore selecting the correct one is vital for performance, previously the relu function has been used as it is the default for sklearn however changing this to identity provides a range of means, ranging from 0.55, a marked improvement, to 0.47. Equally a change to a logistic and tanh activation functions provides a similar range of results, this shows that for this dataset the activation function provides very little change in mean accuracy and therefore improvements must be made elsewhere. This improvement can be seen with the solver functions, where a switch from sgd to ldgfs can provide an improvement to 0.6 with 3000 iterations and a switch to adam can give improvements of up to 0.62 with only 1000 iterations. Finally the teas/train split must be taken into consideration, a change from 0.2 to 0.3 provides a slight decrease in performance, dropping the average mean to 0.56, equally a change to 0.1 also decreases the models accuracy, therefore the originally split of 0.2 is the best. This testing provides the best arrangement of variables for this model on this data set, that arrangement being a hidden layer size of 12, 100 max iterations, the relu activation function, the adam solver, an adaptive learning rate and a test/train split of 0.2 producing an average mean of 0.62, an improvement of 0.28 over the original settings.

Further to the MLP I have implemented a Convolution Neural Network as the second model. This model has some expected similar categories as the MLP, such as the activation functions and test/-train ratios. The CNN also includes a batch size, number of epochs, a softmax function, an optimizer and a loss function. To start with these values are set at using the relu activation function, a 0.2 test/-train split, a 256 batch size, 100 epochs, the rmsprop optimizer, and the sparse categorical crossentropy. This gives a much better accuracy than the original accuracy of the MLP at 0.57. With changes for the activation functions producing no decidable effect on the accuracy with all options producing the same or lower mean accuracy. Next was to test how changing the number of epochs would effect the model, increasing the number from 10 to 1000 did improve the accuracy to 0.65 however the testing accuracy was 0.94, this is indicative of over fitting the model and would make the model less

reliable when interacting with new data, see figure 2. Therefore to make an effective model a balance must be found where the accuracy is high, yet over fitting has not taken place. A value of 250 provides this balance where there is little overfitting and the accuracy has started to plateau, see figure 3. This produced an accuracy of 0.58, which although less than 1000 epochs, is more useful due to the more correct fitting. Attempting to reduce the batch size down to its default of 32 simply shows an increase of overfitting with no accuracy gain showing the 512 batch size is still better, equally the same effect happens when the batch size is doubled to 1024.

For this problem I would recommend the MLP as it consistently provides higher accuracy than the CNN and shows the distinct possibility to be optimized and improved whereas the CNN did not change much when its variables were changed. The low accuracy rates of both these models may be attributed to the problem I raised in the Introduction where the quality measurement is subjective and therefore a clear, generalised pattern may not be possible to find.

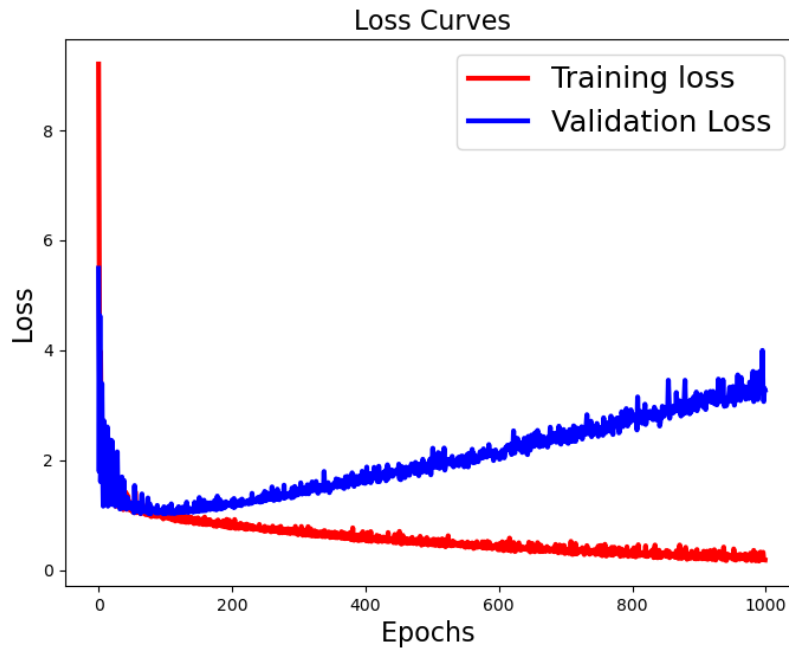


Figure 2: A graph detailing the training and validation loss for the CNN with 1000 epochs.

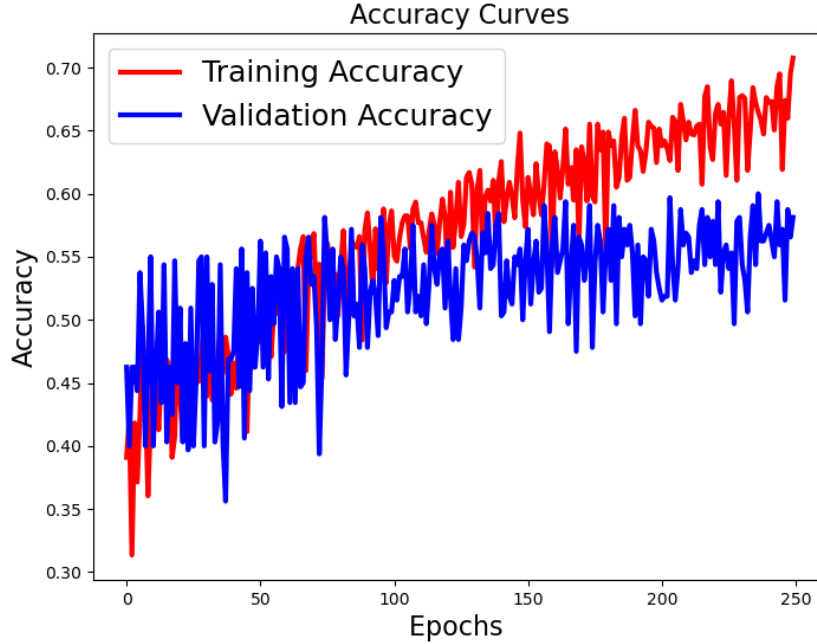


Figure 3: A graph detailing the training and validation accuracy of the CNN with 250 epochs

4 Key Findings and Insights

Having implemented a way to generate randomised wines and pass them to the model I have decided to attempt to predict what the quality of these wines will be and to test if they follow a similar pattern to the previous database and if certain outliers can be found. As can be seen in figure 4 with a small sample of only 100 generated wines, the model greatly favours predicting the absolute average of the classes with 5 taking up the majority of the predictions, what can also be seen is the model will also predict down to 3 and up to 8, however these are in small quantities. This is representative of the original data set as in that set there are no wines that are graded below a 3 or above an 8. From this we can either take away that from a small sample of wines most are very average or that the model is exaggerating its training data and is 'playing it safe'. However when using a much larger data set of 10000 new wines the original trend of favouring a quality of 5, a quality 7 of is relatively high as well, see figure 5. What this might tell us is that even through random generation some variables hold more sway over the

quality and therefore should be more closely watched when creating a wine.

Further to the last two points the problem set out earlier in the introduction has yet been proven further, with that the quality feature for the wine is a subjective measurements in the tester and does not come from a place of true objective fact it makes it hard for the models to truly predict what wine will be good quality as that fact is different fro everybody, although the insite can be used to further improve the wine.

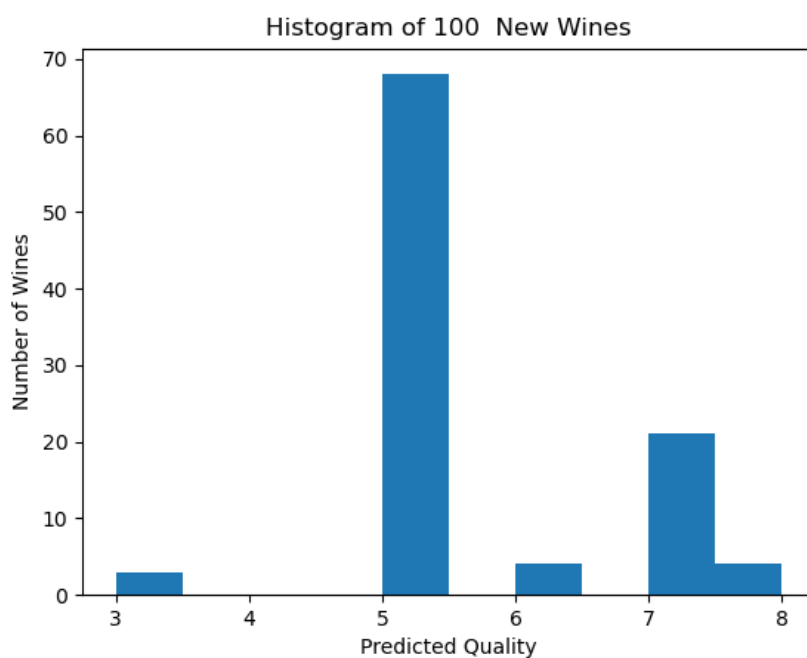


Figure 4: A histogram showing the distribution of predicted qualities for 100 randomly generated wines

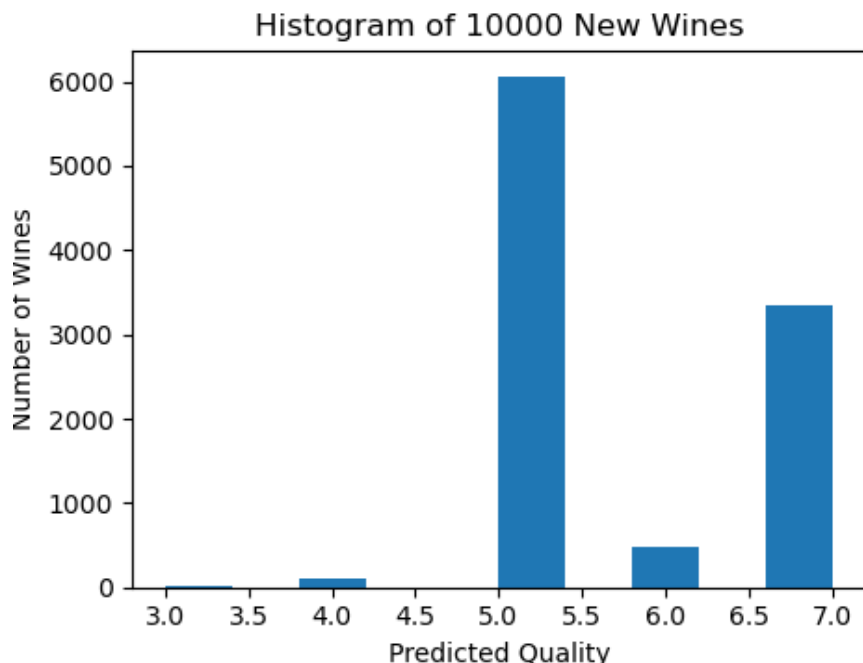


Figure 5: A histogram showing the distribution of predicted qualities for 10000 randomly generated wines

5 Next Steps

The next steps for this analysis is to attempt to produce a more reliable quality measurement perhaps through large scale testing so the model is able to find more of a pattern. Equally a larger data set of wine would help refine the models as 1599 wines may not be enough to distinguish all of the relationships between the compounds. before these issues are fixed I believe different models would not produce much improved results due to the fundamentally subjective nature of the quality. A further issue of the model is that it focuses on the middle ground a lot more and will not produce outputs outside of the set 3-8 range provided by the data set therefore more examples of these sorts of wines, even if there are not many of them would in turn improve the predictive capabilities of the model due to it being able to learn what the qualities of these wines are, potentially being able to recommend near-perfect wines. If this data is collected on one person and the model is used to predict only for that person then the subjective nature of the quality is removed, therefore I can see this being deployed for wine personalisation, where a the model

learns a singular persons tastes and is able to recommend wine types and certain features, although a lot of data would need to be collected on that one person, potentially restricting its effectiveness for a while.

Overall the findings of this project conclude that the MLP model is more consistent with this data than the CNN and that wine quality can be predicted from its chemical composition. However more data and a more objective quality check is required to optimize the models for prediction.

References

- [1] P. Cortez, “Wine quality data set,” 2009. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/Wine+Quality>

6 Appendix 1

```
1 import pandas as pd
2 from sklearn import metrics
3 from sklearn.neural_network import MLPClassifier
4 from sklearn.model_selection import train_test_split
5 from random import randint, uniform
6 import matplotlib.pyplot as plt
7 from tensorflow.keras.models import Sequential
8 from tensorflow.keras.layers import Dense
9 import numpy as np
10
11 # Read in csv file and split by ;
12 wine_df = pd.read_csv('winequality-red.csv', sep=";")
13
14 # Create target dataframe of the quality column
15 target_data = wine_df['quality']
16 del wine_df['quality']
17 # Create dataframe of all other values
18 values = wine_df
19
20 train_features, test_features, train_targets, test_targets =
    train_test_split(values, target_data, test_size=0.2)
21
22
23 def quality_values(target_data, values):
24     '''
25     This function takes in the target data and the values and
26     produces a scatter graph comparing the two,
27     change the y value for different values
28     :param target_data:
29     :param values:
30     '''
31     x = target_data
32     y = values['alcohol']
33
34     # Normalized data to allow for more accurate gradients
35     y_max = y.max()
36     normalized_y = y / y_max
37
38     plt.scatter(x, normalized_y)
39
40     # Create a trend-line
41     z = np.polyfit(x, normalized_y, 1)
42     p = np.poly1d(z)
43     plt.plot(x, p(x), "r--")
44
45     print(z)
46
47     plt.title("Alcohol: Gradient = " + str(round(z[0], 3)))
48     plt.ylabel("Normalised Alcohol")
49     plt.xlabel("Quality")
50     plt.show()
51
52
53 # fixed acidity = small positive correlation = 0.017
```

```

54 # volatile acidity = negative correlation = -0.055
55 # citric acid = positive correlation = 0.055
56 # residual sugar = very small positive correlation, outliers =
    0.002
57 # chlorides = very small negative correlation, outliers = -0.012
58 # free sulfur dioxide = very large negative correlation = -0.009
59 # total sulfur dioxide = very large negative correlation, notable
    outliers = -0.026
60 # density = small negative correlation = 0.0
61 # pH = very small negative correlation, some outliers = -0.003
62 # sulphates = positive correlation, many outliers = 0.026
63 # alcohol = strong positive correlation = 0.042
64
65 def CNN(train_features, test_features, train_targets, test_targets)
    :
    '''
66     Implements a CNN, trains it and returns it
67     :param train_features:
68     :param test_features:
69     :param train_targets:
70     :param test_targets:
71     :return: model: A Tensorflow Sequential Model of the CNN
    '''
72
73     classes_num = 10
74
75     activation = 'relu'
76
77     model = Sequential()
78     model.add(Dense(512, activation=activation, input_shape=(np.
79         prod(train_features.shape[1:]),)))
80     model.add(Dense(512, activation=activation))
81     model.add(Dense(classes_num, activation='softmax'))
82
83     # rmsprop
84     model.compile(optimizer='rmsprop',
85                   loss='sparse_categorical_crossentropy',
86                   metrics=['accuracy'])
87
88     history = model.fit(train_features, train_targets, batch_size
89         =512, epochs=250,
90                       verbose=1, validation_data=(test_features,
91 test_targets))
92
93     [test_loss, test_acc] = model.evaluate(test_features,
94 test_targets)
95     print("Evaluation result on Test Data : Loss = {}, accuracy =
96     {}".format(test_loss, test_acc))
97
98     # Loss and accuracy graphs, uncomment to generate
99
100     # # Plot the Loss Curves
101     # plt.figure(figsize=[8, 6])
102     # plt.plot(history.history['loss'], 'r', linewidth=3.0)
103     # plt.plot(history.history['val_loss'], 'b', linewidth=3.0)
104     # plt.legend(['Training loss', 'Validation Loss'], fontsize=18)
105     # plt.xlabel('Epochs ', fontsize=16)
106     # plt.ylabel('Loss', fontsize=16)

```

```

103 # plt.title('Loss Curves', fontsize=16)
104 #
105 # plt.show()
106 #
107 # # Plot the Accuracy Curves
108 # plt.figure(figsize=[8, 6])
109 # plt.plot(history.history['accuracy'], 'r', linewidth=3.0)
110 # plt.plot(history.history['val_accuracy'], 'b', linewidth=3.0)
111 # plt.legend(['Training Accuracy', 'Validation Accuracy'],
112 #            fontsize=18)
113 # plt.xlabel('Epochs ', fontsize=16)
114 # plt.ylabel('Accuracy', fontsize=16)
115 # plt.title('Accuracy Curves', fontsize=16)
116 #
117 # plt.show()
118
119
120
121 def MLPModel(train_features, test_features, train_targets,
122              test_targets):
123     """
124     Implements a MLP, trains it then returns it
125     :param train_features:
126     :param test_features:
127     :param train_targets:
128     :param test_targets:
129     :return: model: The MLP model
130     """
131     all_means = 0
132     iterations = 50
133     for i in range(iterations):
134         model = MLPClassifier(hidden_layer_sizes=12, max_iter=1000,
135                               activation='relu', solver='adam', verbose=10,
136                               learning_rate='adaptive')
137         model.fit(train_features, train_targets)
138
139         predictions = model.predict(test_features)
140         score = np.round(metrics.accuracy_score(test_targets,
141         predictions), 2)
142         print("Mean accuracy of predictions: " + str(score))
143         all_means += score
144     print("Average means = " + str(all_means / iterations))
145
146     return model
147
148
149
150 def generateWines(values):
151     """
152     This function generates a dataframe of randomly generated wines
153     , without quality measures. The range of each
154     feature is generated by taking the difference between max and
155     min, then subtracting it from the min, to a lower
156     limit of 0 and adding it to the max.
157     :param values:
158     :return: newWines:
159     """

```

```

154 numOfWinesToGenerate = 100
155 newWines = pd.DataFrame()
156
157 minFixedAcidity = values['fixed acidity'].min()
158 maxFixedAcidity = values['fixed acidity'].max()
159 fixedAcidityDifference = maxFixedAcidity - minFixedAcidity
160 fixedAcidityRange = [
161     minFixedAcidity - fixedAcidityDifference if minFixedAcidity
162     - fixedAcidityDifference > 0 else 0,
163     maxFixedAcidity + fixedAcidityDifference]
164
165 minVolatileAcidity = values['volatile acidity'].min()
166 maxVolatileAcidity = values['volatile acidity'].max()
167 volatileAcidityDifference = maxVolatileAcidity -
168 minVolatileAcidity
169 volatileAcidityRange = [
170     minVolatileAcidity - volatileAcidityDifference if
171     minVolatileAcidity - minVolatileAcidity > 0 else 0,
172     maxVolatileAcidity + volatileAcidityDifference]
173
174 minCitricAcid = values['citric acid'].min()
175 maxCitricAcid = values['citric acid'].max()
176 citricAcidDifference = maxCitricAcid - minCitricAcid
177 citricAcidRange = [minCitricAcid - citricAcidDifference if
178     minCitricAcid - citricAcidDifference > 0 else 0,
179     maxCitricAcid + citricAcidDifference]
180
181 minResidualSugar = values['residual sugar'].min()
182 maxResidualSugar = values['residual sugar'].max()
183 residualSugarDifference = maxResidualSugar - minResidualSugar
184 residualSugarRange = [
185     minResidualSugar - residualSugarDifference if
186     minResidualSugar - residualSugarDifference > 0 else 0,
187     maxResidualSugar + residualSugarDifference]
188
189 minChlorides = values['chlorides'].min()
190 maxChlorides = values['chlorides'].max()
191 chloridesDifference = maxChlorides - minChlorides
192 chloridesRange = [minChlorides - chloridesDifference if
193     minChlorides - chloridesDifference > 0 else 0,
194     maxChlorides + chloridesDifference]
195
196 minFreeSulphurDioxide = values['free sulfur dioxide'].min()
197 maxFreeSulphurDioxide = values['free sulfur dioxide'].max()
198 freeSulphurDioxideDifference = maxFreeSulphurDioxide -
199 minFreeSulphurDioxide
200 freeSulphurDioxideRange = [
201     minFreeSulphurDioxide - freeSulphurDioxideDifference if
202     minFreeSulphurDioxide - freeSulphurDioxideDifference > 0
203     else 0,
204     maxFreeSulphurDioxide + freeSulphurDioxideDifference]
205
206 minTotalSulphurDioxide = values['total sulfur dioxide'].min()
207 maxTotalSulphurDioxide = values['total sulfur dioxide'].max()
208 totalSulphurDioxideDifference = maxTotalSulphurDioxide -
209 minTotalSulphurDioxide
210 totalSulphurDioxideRange = [

```



```

202     minTotalSulphurDioxide - totalSulphurDioxideDifference if
203     minTotalSulphurDioxide - totalSulphurDioxideDifference > 0
    else 0,
204     maxTotalSulphurDioxide + totalSulphurDioxideDifference]
205
206     minDensity = values['density'].min()
207     maxDensity = values['density'].max()
208     densityDifference = maxDensity - minDensity
209     densityRange = [minDensity - densityDifference if minDensity -
    densityDifference > 0 else 0,
210                     maxDensity + densityDifference]
211
212     minpH = values['pH'].min()
213     maxpH = values['pH'].max()
214     pHDifference = maxpH - minpH
215     pHRange = [minpH - pHDifference if minpH - pHDifference > 0
    else 0, maxpH + pHDifference]
216
217     minSulphates = values['sulphates'].min()
218     maxSulphates = values['sulphates'].max()
219     sulphatesDifference = maxSulphates - minSulphates
220     sulphatesRange = [minSulphates - sulphatesDifference if
    minSulphates - sulphatesDifference > 0 else 0,
221                       maxSulphates + sulphatesDifference]
222
223     minAlcohol = values['alcohol'].min()
224     maxAlcohol = values['alcohol'].max()
225     alcoholDifference = maxAlcohol - minAlcohol
226     alcoholRange = [minAlcohol - alcoholDifference if minAlcohol -
    alcoholDifference > 0 else 0,
227                     maxAlcohol + alcoholDifference]
228
229     for i in range(numOfWinesToGenerate):
230         newWines = newWines.append([[uniform(fixedAcidityRange[0],
    fixedAcidityRange[1]),
231                                     uniform(volatileAcidityRange
    [0], volatileAcidityRange[1]),
232                                     uniform(citricAcidRange[0],
    citricAcidRange[1]),
233                                     uniform(residualSugarRange[0],
    residualSugarRange[1]),
234                                     uniform(chloridesRange[0],
    chloridesRange[1]),
235                                     randint(
    freeSulphurDioxideRange[0], freeSulphurDioxideRange[1]),
236                                     randint(
    totalSulphurDioxideRange[0], totalSulphurDioxideRange[1]),
237                                     uniform(densityRange[0],
    densityRange[1]),
238                                     uniform(pHRange[0], pHRange
    [1]),
239                                     uniform(sulphatesRange[0],
    sulphatesRange[1]),
240                                     (uniform(alcoholRange[0],
    alcoholRange[1]))]], ignore_index=True)
241     return newWines
242

```

```

243
244 def predictQuality(newWines, MLPModel, CNNModel):
245     """
246     A simple function that just decides which model to use, it
247     returns an array of the predictions
248     :param newWines:
249     :param MLPModel:
250     :param CNNModel:
251     :return: predictions: Array of the predictions
252     """
253     if MLPModel is not None:
254         predictions = MLPModel.predict(newWines)
255     else:
256         predictions = CNNModel.predict_on_batch(newWines)
257
258     return predictions
259
260 mode = 0
261 if mode == 0:
262     prediction = predictQuality(generateWines(values),
263                                MLPModel(train_features,
264                                test_features, train_targets, test_targets),
265                                None)
266 else:
267     prediction = predictQuality(generateWines(values),
268                                None,
269                                CNN(train_features, test_features,
270                                train_targets, test_targets))
271
272 # A histogram of the predictions generated by the model
273 plt.hist(prediction, bins=10, histtype='stepfilled')
274 plt.title('Histogram of 100 New Wines')
275 plt.xlabel('Predicted Quality')
276 plt.ylabel('Number of Wines')
277 plt.show()

```

7 Appendix 2

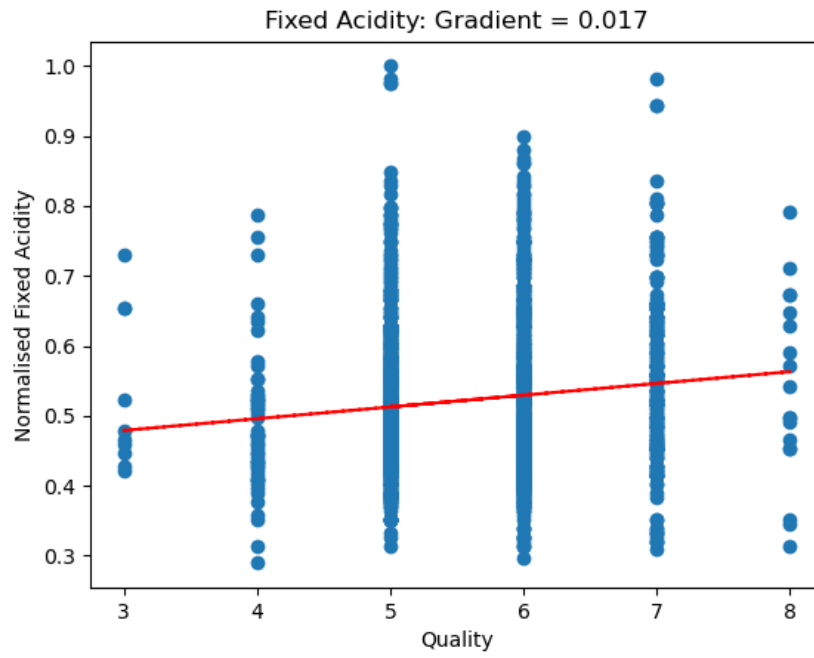


Figure 6: A graph showing the relationship between fixed acidity and quality

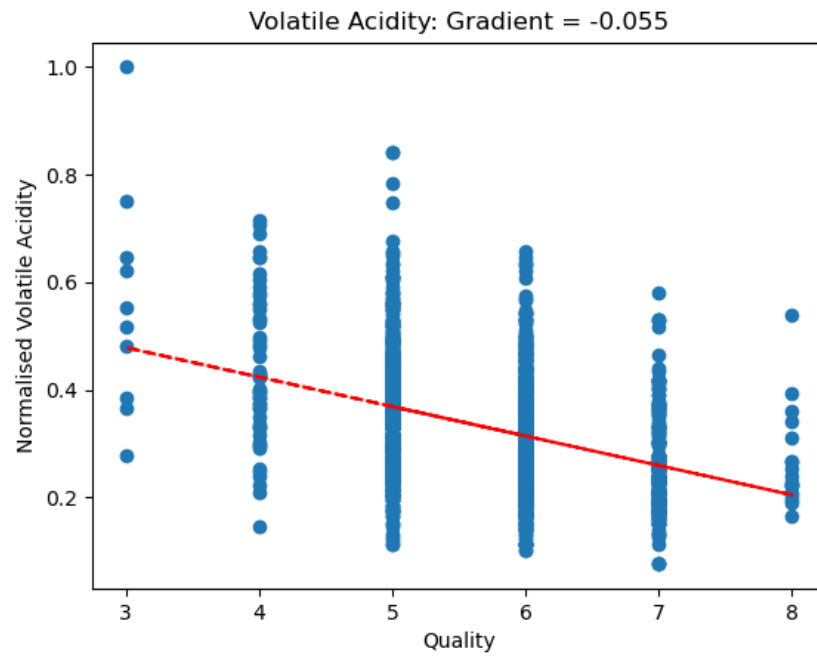


Figure 7: A graph showing the relationship between volatile acidity and quality

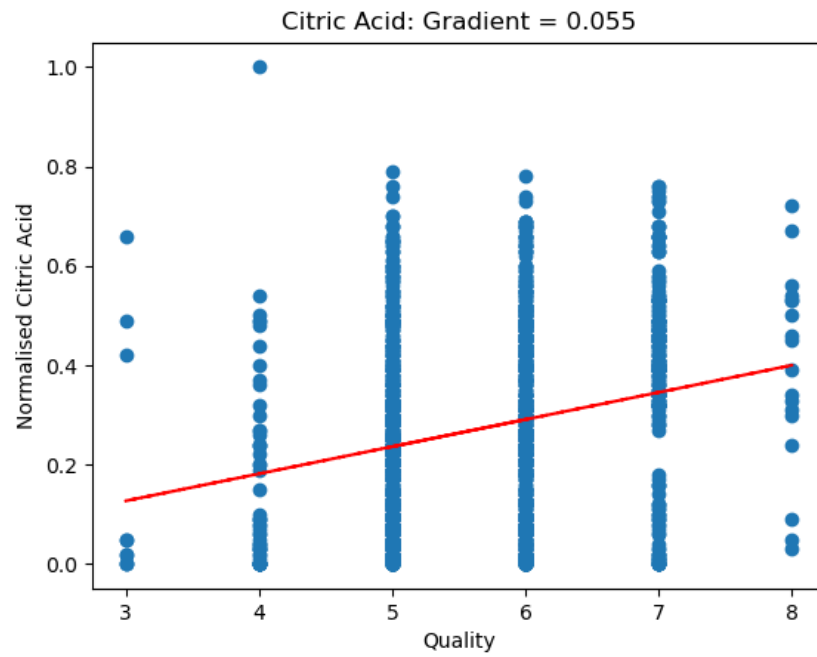


Figure 8: A graph showing the relationship between citric acid and quality

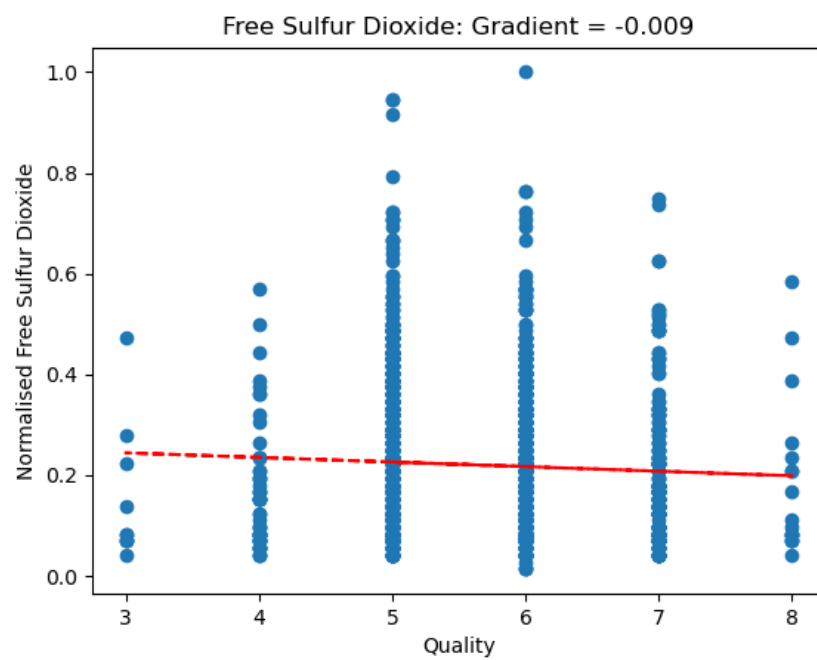


Figure 9: A graph showing the relationship between free sulphur dioxide and quality

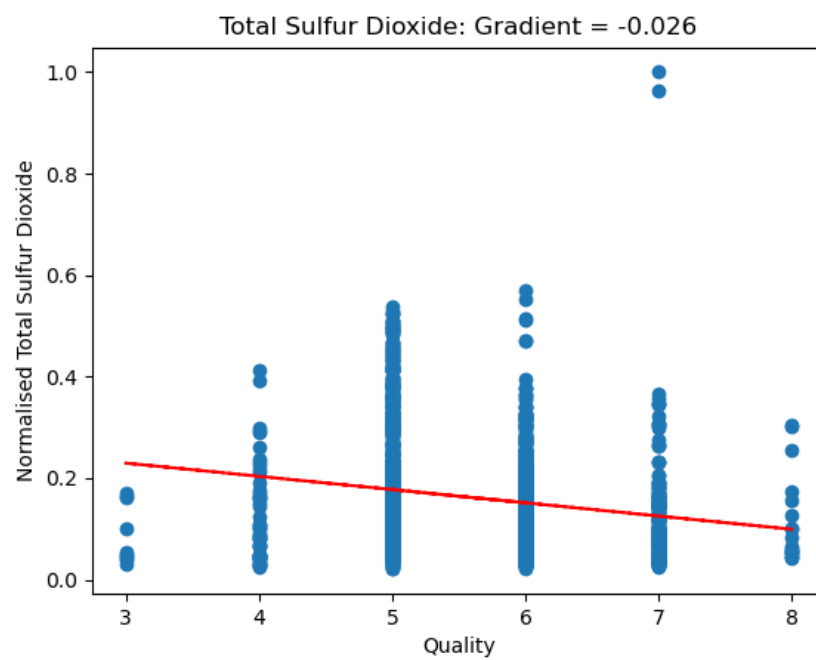


Figure 10: A graph showing the relationship between total sulphur dioxide and quality

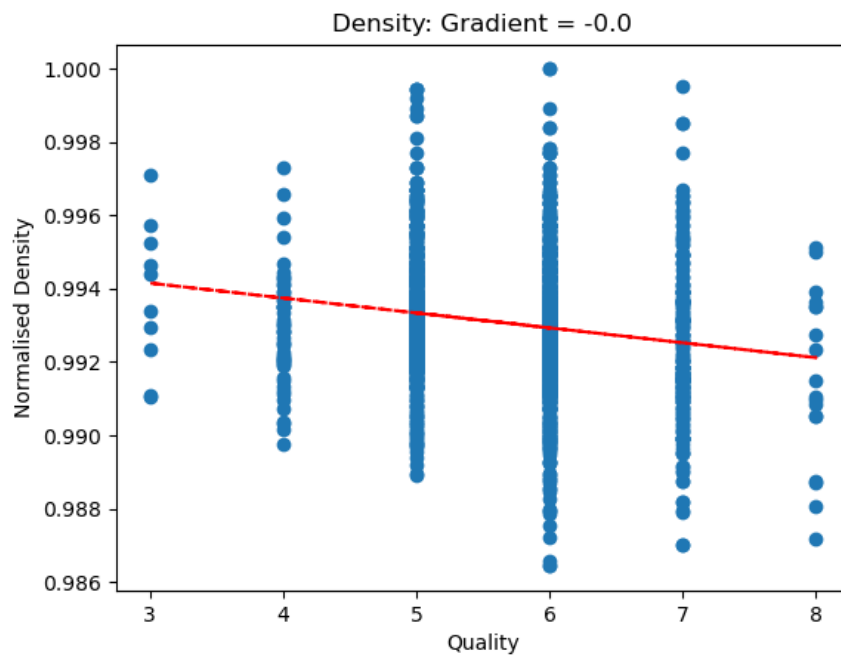


Figure 11: A graph showing the relationship between density and quality

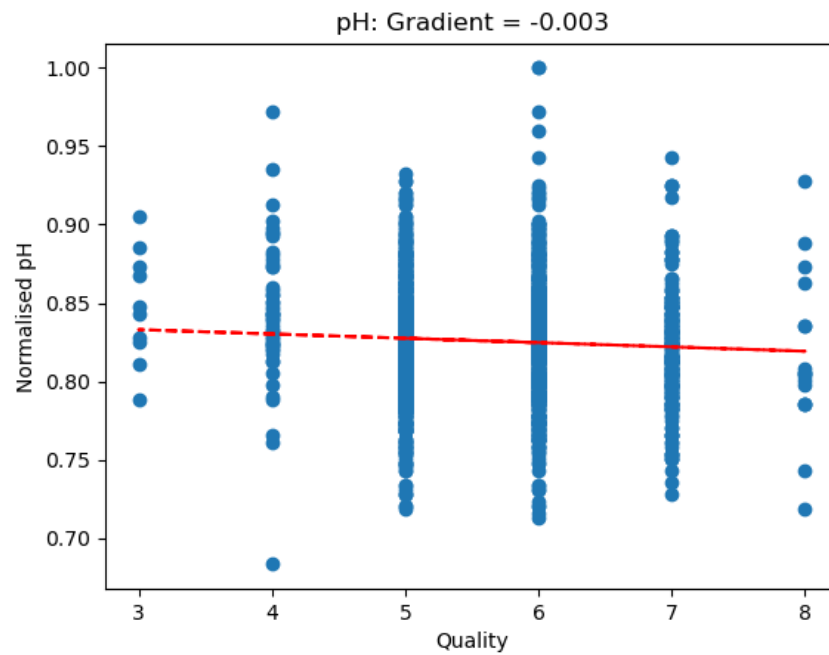


Figure 12: A graph showing the relationship between pH and quality

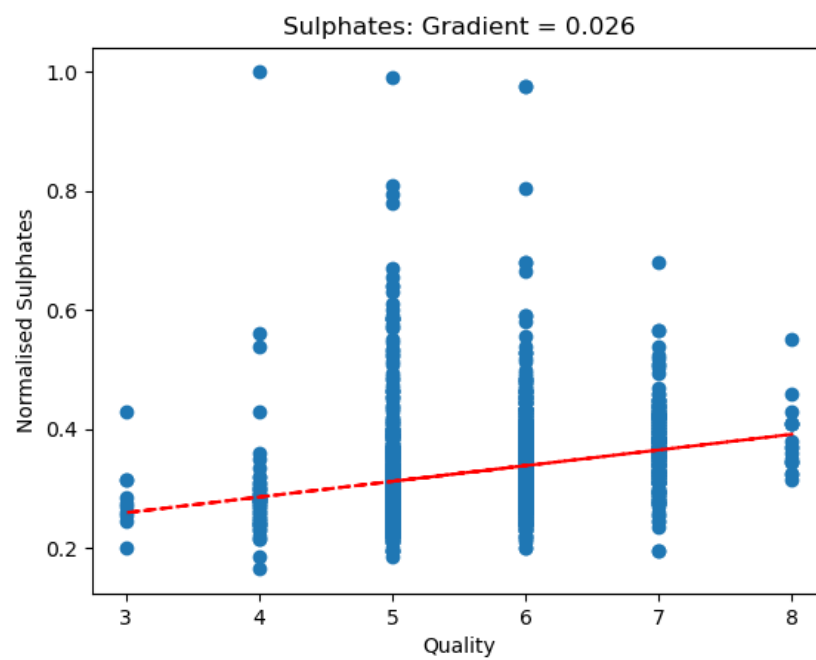


Figure 13: A graph showing the relationship between sulphates and quality

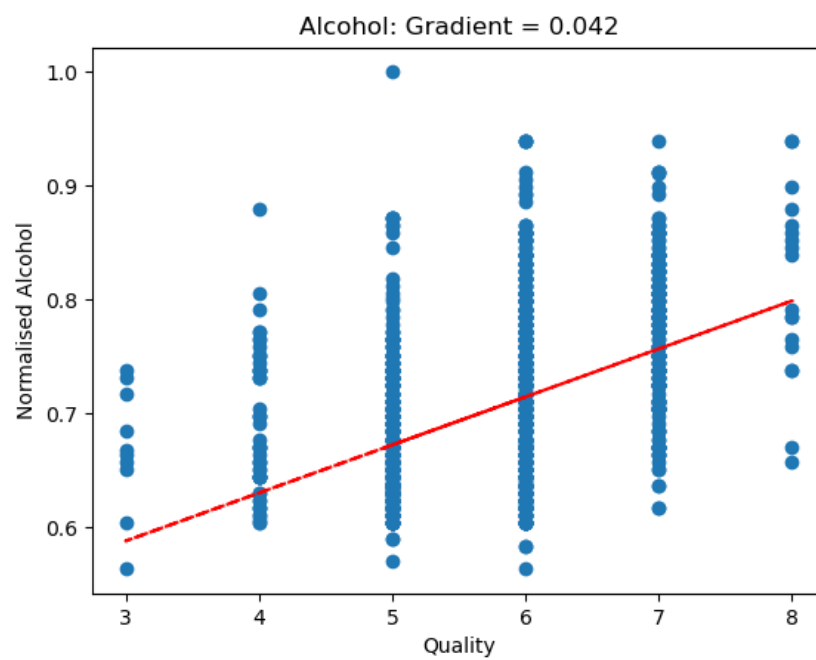


Figure 14: A graph showing the relationship between alcohol and quality