

```

1 %%Resetting the environment
2 %clear the command window and any variables, and close all figure windows
3 clc;
4 clear;
5 close all;
6
7
8 %%Loading the dataset
9 %link for the dataset: https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009?resource=download
10 %Using readtable function because our dataset contains not only numeric
11 %data but also numeric and text columns
12 redwine = readtable('winequality-red.csv');
13 % Display the first few rows of the dataset to understand its structure
14 head(redwine);
15 summary(redwine);
16
17 %%Check for any missing values using the ismissing function
18 %reference: https://uk.mathworks.com/help/matlab/ref/ismissing.html
19 missingValues = sum(ismissing(redwine));
20 missingValues;
21 %The output shows no missing values for all the variables so we are fine
22
23 %%Now we will remove any duplicate rows and then we will
24 % display the number of duplicate rows removed so we can tell how many rows
25 % have been removed
26 %reference: https://www.mathworks.com/matlabcentral/answers/275640-deleting-repeated-rows-in-a-matrix
27 Before = height(redwine);
28 redwine = unique(redwine, "rows");
29 After = height(redwine);
30 fprintf('Number of duplicate rows removed: %d\n', Before - After);
31
32 %%Now, since our problem has to do with the prediction whether the wine is
33 %of Low, Medium or High quality, the next step is to convert the quality column
34 %into a binary classification, where 3,4 is Low, 5,6 is Medium and 7,8
35 %is High quality and create a new column named QualityLabel
36 q = redwine.quality;
37
38 Quality = strings(size(q));
39 Quality(q <= 4) = "Low";
40 Quality(q >= 5 & q <= 6) = "Medium";
41 Quality(q >= 7) = "High";
42 redwine.QualityLabel = categorical(Quality, ["Low", "Medium", "High"]);
43
44 disp("Distribution of Class (Low, Medium, High):");
45 summary(redwine.QualityLabel);
46
47 %%Next step is to select only the predictor variables (the 11
48 %%physicochemical input features) and define the categorical target
49 %%variable
50 predictors = setdiff(redwine.Properties.VariableNames, {'quality', 'QualityLabel'});%Here
we used the setdiff function
51 %in order to keep all the columns except quality and QualityLabel which are
52 %the target variables
53 Xtable = redwine(:, predictors);
54 Y = redwine.QualityLabel
55
56 %%Now we will standardise our data using the z-score
57 X = table2array(redwine(:, predictors));
58 Xz = zscore(X);

```

```

59
60 %%Now we will visualise the predictors with the standardised values using a
61 %%boxplot
62 %reference: https://uk.mathworks.com/help/stats/boxplot.html
63 boxplot(Xz, 'Labels', predictors);
64 title('Boxplot of Physicochemical Features');
65 xlabel('Features');
66 ylabel('Values(standardised)');
67 xtickangle(45);
68 grid on;
69
70 %%Now we will create a heatmap for the correlation between predictors
71 %reference: https://uk.mathworks.com/help/matlab/ref/heatmap.html
72 %First, we create a correlation matrix
73 correlationMatrix = corr(redwine(:, predictors));
74 figure;
75 heatmap = heatmap(predictors, predictors, correlationMatrix);
76 heatmap.Colormap = hot;
77 title('Correlation Matrix of predictors');
78
79 %%Now we will remove any outliers. We will do this for Feature-level and
80 %%Instance-level and compare which method has the strongest suspects of
81 %%outliers.
82
83 %For Feature-level we will use the IQR method
84 %reference: https://uk.mathworks.com/help/matlab/ref/iqr.html
85 %reference: https://uk.mathworks.com/help/matlab/ref/isoutlier.html
86 n = size(Xz, 1);
87 p = size(Xz, 2);
88
89 featureOutlierMask = false(n, p);
90
91 for j=1:p
92     col = Xz(:, j);
93     Q1 = quantile(col, 0.25);
94     Q3 = quantile(col, 0.75);
95     IQR = Q3 - Q1;
96     lowerBound = Q1 - 1.5 * IQR;
97     upperBound = Q3 + 1.5 * IQR;
98     featureOutlierMask(:, j) = col < lowerBound | col > upperBound;
99 end
100
101 nFeatureOutliers = sum(featureOutlierMask, 2);
102 redwine.nFeatureOutliers = nFeatureOutliers;
103
104 %Now we set a strong rule for outlying at least 3 features
105 featureLevelFlag = (nFeatureOutliers >= 3);
106
107 fprintf("Feature-level flagged rows: %d\n", sum(featureLevelFlag));
108 %So we have 25 flagged rows in Feature-level
109
110 %Now we move on to Instance-level. We will use the Mahalanobis distance
111 %method. We will compute squared the distance of each row to the
112 %distribution
113 %reference: ChatGPT: Help me detect any outliers for Instance-level and the
114 %find the strongest suspects comparing with the Feature-level outlier
115 %detection.
116 D2 = mahal(Xz, Xz);
117 redwine.MahalanobisD2 = D2;
118
119 %Now, we flag the top 3% as instance-level outliers

```

```

120 contamination = 0.03;
121 cutoff = quantile(D2, 1-contamination);
122 instanceLevelFlag = (D2>=cutoff);
123
124 fprintf("Instance-level flagged rows: %d\n", sum(instanceLevelFlag));
125
126 %So we have 41 flagged rows in Instance-level
127
128 % Now we will combine the feature-level and instance-level flags to identify rows
129 % flagged as outliers
130 outliers = featureLevelFlag & instanceLevelFlag;
131 suspects = find(outliers);
132
133 fprintf("Strongest outlier suspects: %d\n", numel(suspects));
134
135 %Now we will display the results in a table
136 %reference: ChatGPT: Create a table for the suspects results
137 TableOfSuspects = redwine(suspects, :);
138 disp("Outlier suspect row indices:");
139 disp(suspects);
140 disp("Preview of suspected outlier rows:");
141 disp(TableOfSuspects(:, [predictors, {'quality',
142 'QualityLabel','nFeatureOutliers','MahalanobisD2'}]));
143
144 %%Now, before and after removing the outliers, we will perform PCA 2D
145 %%visualisation in order to see the differences
146 %reference: https://uk.mathworks.com/help/stats/pca.html
147 %reference: https://uk.mathworks.com/help/stats/gscatter.html
148 [coeff, score] = pca(Xz);
149
150 figure;
151 gscatter(score(:,1), score(:,2), redwine.QualityLabel);
152 hold on;
153 plot(score(outliers,1), score(outliers,2), 'ko', 'MarkerSize',8,'LineWidth',1.5);
154 hold off;
155 title("PCA visualisation for outliers suspects (before removing them)");
156 xlabel('PC1');
157 ylabel('PC2');
158 grid on;
159
160 % Now we will remove the outliers from the dataset and convert a new
161 % cleaned dataset, X table for the predictors, Xz score and Y table for the
162 % QualityLabel
163 redwine_clean = redwine(~outliers, :);
164 X_clean = table2array(redwine_clean(:, predictors));
165 Xz_clean = zscore(X_clean);
166 Y_clean = redwine_clean.QualityLabel;
167
168 fprintf("Rows before removing the outliers: %d\n", height(redwine));
169 fprintf("Rows after removing the outliers: %d\n", height(redwine_clean));
170
171 %And now we will perform PCA visualisation again, after removing the
172 %outliers
173 [~, score2] = pca(Xz_clean);
174
175 figure;
176 gscatter(score2(:,1), score2(:,2), Y_clean);
177 title("PCA visualisation (after removing the outliers)");
178 xlabel('PC1');
179 ylabel('PC2');
180 grid on;

```

```

179
180 %%This concludes the pre processing, cleaning and visualising the data. Now
181 %%we can move on to the machine learning algorithms
182
183 %We will train/test the two algorithms with the same techniques for fair
184 %comparison
185 %Firstly, we set a random number generator for reproducibility
186 rng(1)
187
188 %%Now we will split the dataset into training and test (70% train and 30%
189 %%test)
190 %reference: https://www.mathworks.com/matlabcentral/answers/377839-split-training-data-and-testing-data?s\_tid=srchtitle
191 %First we create a partition
192 cv = cvpartition(Y_clean, 'Holdout', 0.3);
193 trainidx=cv.training;
194 testidx=cv.test;
195
196 %Now we seperate to training and test data
197 Xtrain = Xz_clean(trainidx,:);
198 Ytrain = Y_clean(trainidx,:);
199
200 Xtest = Xz_clean(testidx,:);
201 Ytest = Y_clean(testidx,:);
202
203 %Here we order the classes
204 classOrder = categories(Y_clean);
205
206 %%Next step is to train the Naive Bayes model
207 %reference: https://uk.mathworks.com/help/stats/fitcnb.html
208 %Also, we follow the lecture notes and the tutorial from week 4 of the
209 %module : https://moodle4.city.ac.uk/course/section.php?id=243698
210 NaiveBayes = fitcnb(Xtrain, Ytrain,'DistributionNames','normal');
211
212 %Prediction for our Naive Bayes trained model
213 [YpredictedNaiveBayes, posteriorNaiveBayes] = predict(NaiveBayes, Xtest);
214
215 %%Now we will evaluate the Naive Bayes classifier
216 %Starting with the confusion matrix
217 %reference: https://uk.mathworks.com/help/stats/confusionchart.html
218 confusionmatrix = confusionmat(Ytest, YpredictedNaiveBayes, 'Order', classOrder);
219 figure;
220 confusionchart(Ytest, YpredictedNaiveBayes);
221 title('Confusion Matrix of Naive Bayes');
222
223 %Now we use 5-fold cross-validation to give a more reliable estimate of
224 %performance of our data and to detect overfitting
225 %reference: https://uk.mathworks.com/help/stats/crossval.html
226 crossvalidationModel = crossval(NaiveBayes, 'KFold', 5);
227 crossvalidationLoss = kfoldLoss(crossvalidationModel);
228 fprintf('Cross-validation loss for the Naive Bayes classifier: %.3f\n',
crossvalidationLoss);
229
230 %Now we will calculate the metrics for the Naive Bayes (f1 score,
231 %precision and recall)
232 %We will first find the true positive, true negative, false positive and
233 %false negative values
234 %reference: https://uk.mathworks.com/matlabcentral/answers/1703920-how-can-i-calculate-the-false-negatives-false-positives-and-the-correctly-classified-for-my-seizure
235 classOrder = categories(Y_clean);
236

```

```

237 TPnaivebayes = diag(confusionmatrix);
238 FPnaivebayes = sum(confusionmatrix,1)' - TPnaivebayes;
239 FNnaivebayes = sum(confusionmatrix,2) - TPnaivebayes;
240 TNnaivebayes = sum(confusionmatrix(:)) - TPnaivebayes - FPnaivebayes - FNnaivebayes;
241
242 %Now we calculate the metrics
243 %reference: https://medium.com/analytics-vidhya/confusion-matrix-accuracy-precision-recall-f1-score-ade299cf63cd
244
245 %Now we calculate the precision
246 precisionNaiveBayes = TPnaivebayes ./ max(TPnaivebayes+FPnaivebayes,1);
247
248 %Now we calculate the recall
249 recallNaiveBayes = TPnaivebayes ./ max(TPnaivebayes + FNnaivebayes,1);
250
251 %Now we calculate the f1 score
252 f1scoreNaiveBayes = 2*(precisionNaiveBayes .* recallNaiveBayes) ./ max(precisionNaiveBayes+recallNaiveBayes, eps);
253
254 fprintf('Per class precision for Naive Bayes:'); disp(precisionNaiveBayes);
255 fprintf('Per class recall for Naive Bayes:'); disp(recallNaiveBayes);
256 fprintf('Per class f1 score for Naive Bayes:'); disp(f1scoreNaiveBayes);
257
258 %Now, because in multiclass problems we don't have single metrics we will
259 %compute the macro-metrics(average of the per class metrics for equal weight across
260 %classes) and the
261 %weighted-metrics(average of the per class metrics weighted by support)
262 %reference: chatGPT: help me compute the macro and weighted metrics
263 supportNaiveBayes = sum(confusionmatrix,2);
264
265 macroPrecisionNaiveBayes = mean(precisionNaiveBayes);
266 macroRecallNaiveBayes = mean(recallNaiveBayes);
267 macrof1scoreNaiveBayes = mean(f1scoreNaiveBayes);
268
269 weightedPrecisionNaiveBayes = sum(precisionNaiveBayes .* supportNaiveBayes) /
270 sum(supportNaiveBayes);
271 weightedRecallNaiveBayes = sum(recallNaiveBayes .* supportNaiveBayes) /
272 sum(supportNaiveBayes);
273 weightedf1scoreNaiveBayes = sum(f1scoreNaiveBayes .* supportNaiveBayes) /
274 sum(supportNaiveBayes);
275
276 %Now we creat a table to represent the metrics per class
277 perClassNaiveBayes = table(classOrder, TPnaivebayes, TNnaivebayes,
278 FPnaivebayes, FNnaivebayes, precisionNaiveBayes, recallNaiveBayes,
279 f1scoreNaiveBayes, supportNaiveBayes, 'VariableNames',{'Class','TP',
280 'TN','FP','FN','Precision','Recall','f1score', 'Support'});
281 disp('Naive Bayes: Metrics per class');
282 disp(perClassNaiveBayes);
283
284 %%And we are done with the Naive Bayes
285
286 %%Now we move on to the Decision Tree model. We will use same training data
287 %%for a fair comparison
288 %reference: https://uk.mathworks.com/help/stats/fitctree.html? s\_tid=srchtitle\_support\_results\_1\_train+decision+tree
289 DecisionTree = fitctree(Xtrain, Ytrain, 'SplitCriterion','gdi','MinLeafSize',10);
290
291 %Prediction for our Decision Tree trained model
292 [YpredictedDecisionTree, DecisionTreeScores] = predict(DecisionTree, Xtest);
293
294 %%Visualization of the decision tree

```

```

288 %reference: https://
uk.mathworks.com/support/search.html/answers/1567223-visualization-of-decision-tree.html?
fq%5B%5D=asset_type_name:answer&fq%5B%5D=category:stats/data-import-and-export-1&page=1
289 figure;
290 view(DecisionTree, 'Mode', 'graph');
291 title('Decision Tree');
292
293 %%Now we will evaluate the Decision Tree classifier
294 %Starting with the confusion matrix
295 %reference: https://uk.mathworks.com/help/stats/confusionchart.html
296 confusionmatrixTree = confusionmat(Ytest, YpredictedDecisionTree, 'Order', classOrder);
297 figure;
298 confusionchart(Ytest, YpredictedDecisionTree);
299 title('Confusion Matrix of Decision Tree');
300
301 % Now we will evaluate the performance of the Decision Tree classifier using cross-validation
302 crossvalidationModelTree = crossval(DecisionTree, 'KFold', 5);
303 crossvalidationLossTree = kfoldLoss(crossvalidationModelTree);
304 fprintf('Cross-validation loss for the Decision Tree classifier: %.3f\n',
crossvalidationLossTree);
305
306 %Now we will calculate the metrics for the Decision Tree (f1 score,
307 %precision and recall)
308 %We will first find the true positive, true negative, false positive and
309 %false negative values
310 classOrder = categories(Y_clean);
311
312 TPTree = diag(confusionmatrixTree);
313 FPTree = sum(confusionmatrixTree,1)' - TPTree;
314 FNTree = sum(confusionmatrixTree,2) - TPTree;
315 TNTree= sum(confusionmatrixTree(:)) - TPTree - FPTree - FNTree;
316
317 %Now we calculate the metrics
318 %reference: https://medium.com/analytics-vidhya/confusion-matrix-accuracy-precision-
recall-f1-score-ade299cf63cd
319
320 %Now we calculate the precision
321 precisionTree = TPTree ./ max(TPTree+FPTree,1);
322
323 %Now we calculate the recall
324 recallTree = TPTree./ max(TPTree + FNTree,1);
325
326 %Now we calculate the f1 score
327 f1scoreTree= 2*(precisionTree .* recallTree) ./ max(precisionTree+recallTree, eps);
328
329 fprintf('precision for Decision Tree:'); disp(precisionTree);
330 fprintf('recall for Decision Tree:'); disp(recallTree);
331 fprintf('f1 score for Decision Tree:'); disp(f1scoreTree);
332
333 %Now, because in multiclass problems we don't have single metrics we will
334 %compute the macro-metrics(average of the per class metrics for equal weight across
classes) and the
335 %weighted-metrics(average of the per class metrics weighted by support)
336 %reference: chatGPT: help me compute the macro and weighted metrics
337 supportTree = sum(confusionmatrixTree,2);
338
339 macroPrecisionTree = mean(precisionTree);
340 macroRecallTree = mean(recallTree);
341 macrof1scoreTree = mean(f1scoreTree);
342

```

```
343 weightedPrecisionTree = sum(precisionTree .* supportTree) / sum(supportTree);
344 weightedRecallTree = sum(recallTree .* supportTree) / sum(supportTree);
345 weightedf1scoreTree = sum(f1scoreTree .* supportTree) / sum(supportTree);
346
347 %Now we creat a table to represent the metrics per class
348 perClassTree = table(classOrder, TPTree, TNTree, FPTree, FNTree,
precisionTree, recallTree, f1scoreTree, supportTree, 'VariableNames',{'Class','TP',
'TN','FP','FN','Precision','Recall','f1score','Support'});
349 disp('Decision Tree: Metrics per class');
350 disp(perClassTree);
351
352 % Now we will visualize the performance metrics for both classifiers
353 % Now we will plot the performance metrics for both classifiers
354 figure;
355 bar([macroPrecisionTree, macroRecallTree, macrof1scoreTree; weightedPrecisionTree,
weightedRecallTree, weightedf1scoreTree]);
356 set(gca, 'XTickLabel', {'Macro', 'Weighted'});
357 ylabel('Scores');
358 title('Performance Metrics Comparison');
359 legend('Precision', 'Recall', 'F1 Score');
360
361 %And this concludes the code
362
```