

FIT3152 Data analytics

Assignment 2

Student Name: Kuah Jia Chen

Student ID: 32286988

Overview

Please refer to the Appendix for all the R code I used for the analysis in this report. There are a total of 11 pages in this report. (i.e., the number of pages for Appendix is not counted here).

Question 1)

Answer:

The proportion of days when it is warmer than the previous day is 54.6%, whereas the proportion of days when it is cooler than the previous day is 44.7%. We can see that the proportion of days when it is warmer than the previous day is higher than those where it is cooler by roughly 10%, hence more days are warmer than the previous day. There are 13 rows with a missing value for the WarmerTomorrow attribute. (i.e., about 0.7%)

Furthermore, I used the summary function as well as the sd function to get the value of minimum, 1st quartile, median, mean, 3rd quartile, the number of NAs, and the standard deviation to get the description for real-valued attributes. Those categorical attributes are not included.

```
> summary(airquality[1:19,15:19,22:23]) # use summary function only for those columns that are real-valued
      MinTemp      MaxTemp      Rainfall      Evaporation      Sunshine      windgustspeed      windgustmax      windgust15m      Humidity9am
 Min: -7.4   Min: 17.4   Min: 0.00   Min: 0.0   Min: 0.0   Min: 1.00   Min: 1.00   Min: 1.00   Min: 1.00
 1st Qu: 6.8   1st Qu: 24.2   1st Qu: 0.00   1st Qu: 0.0   1st Qu: 0.0   1st Qu: 11.0   1st Qu: 11.0   1st Qu: 11.0   1st Qu: 22.0
 Median: 11.8   Median: 29.8   Median: 0.00   Median: 0.0   Median: 0.0   Median: 11.0   Median: 11.0   Median: 11.0   Median: 72.0
 Mean: 12.7   Mean: 27.8   Mean: 1.38   Mean: 0.0   Mean: 0.0   Mean: 12.3   Mean: 12.3   Mean: 12.3   Mean: 68.9
 3rd Qu: 17.2   3rd Qu: 34.2   3rd Qu: 0.00   3rd Qu: 0.0   3rd Qu: 0.0   3rd Qu: 13.0   3rd Qu: 13.0   3rd Qu: 13.0   3rd Qu: 82.0
 Max: 18.8   Max: 45.7   Max: 115.4   Max: 10.0   Max: 11.0   Max: 13.0   Max: 13.0   Max: 13.0   Max: 100.0
 NA's: 19   NA's: 19   NA's: 19   NA's: 1270   NA's: 1270   NA's: 1270   NA's: 1270   NA's: 1270   NA's: 19

      Humidity3pm      Pressure9am      Pressure3pm      Temp9am      Temp3pm
 Min: 1.00   Min: 1.000   Min: 1.000   Min: -4.4   Min: -2.7
 1st Qu: 10.0   1st Qu: 1000   1st Qu: 1000   1st Qu: 12.3   1st Qu: 10.9
 Median: 13.0   Median: 1000   Median: 1000   Median: 12.0   Median: 12.9
 Mean: 12.6   Mean: 1100.0   Mean: 1100.0   Mean: 12.6   Mean: 12.1
 3rd Qu: 16.0   3rd Qu: 1000   3rd Qu: 1000   3rd Qu: 12.8   3rd Qu: 12.3
 Max: 180.0   Max: 1100.0   Max: 1100.0   Max: 15.7   Max: 13.8
 NA's: 137   NA's: 1000   NA's: 1000   NA's: 146   NA's: 146
```

Figure 1 shows the description of the predictor variables with real-valued attributes

```
> standard_deviation_dataframe
      variable standard deviation
1      MinTemp           5.66
2      MaxTemp           7.43
3      Rainfall           9.64
4      Evaporation        11.8
5      Sunshine          17.4
6      windgustspeed       9.43
7      windgustmax        13.0
8      windgust15m        13.0
9      Humidity9am        18.4
10     Humidity3pm        20.1
11     Pressure9am         6.60
12     Pressure3pm         6.60
13     Temp9am            0.93
14     Temp3pm            0.93
```

Figure 2 shows the standard deviation for predictor variables for real-valued attributes

Figure 1 and Figure 2 show the description of the real-valued attributes. Since the standard deviation is not provided when using the summary function, therefore I calculate it and store the results in a separate data frame.

In my opinion, the most noteworthy thing in the data for those columns is that there are so many NA values in each column. Especially for Evaporation and Sunshine, they have more than a thousand NA values as shown in Figure 1. Besides that, we also can see that Humidity3pm has the highest value of standard deviation, which is around 20, followed by Humidity9am. Furthermore, since the units for each variable are not consistent and unrelated to each other, therefore their mean, median, and so on as shown in Figure 1 are not comparable. Nevertheless, we can see that for MinTemp, MaxTemp, Temp9am and Temp3pm, their minimum values are negative. Moreover, it is also interesting to note that the difference between the minimum value and maximum value for Rainfall is quite significant (i.e., about 115).

Additionally, there are some attributes that I need to consider omitting from my analysis, which are Day, Month, and Year. The reason is that these three attributes only represent the date of the observation, it does not provide any useful information to predict whether tomorrow will be warmer than today. Hence, I think that including them will cause noise to the data. Initially, I was thinking of omitting Location and the attributes that related to the

direction of the wind as well. However, I had done a bit of research online and realized that these attributes do contribute to the temperature, thus I eventually decided to keep them. In conclusion, I only omitted Day, Month, and Year from my analysis. (i.e., I will remove these attributes in Part 2)

Question 2)

Answer:

Before I perform any model fitting, I need to do some pre-processing on the data set to make it suitable to fit the model. Firstly, I will need to ignore Day, Month, and Year from my analysis due to the reason stated in Part 1. Therefore, I will create a new data frame called “new.WAUS” and it is the same as “WAUS”, except that it does not have the attributes Day, Month, and Year.

After that, I need to convert the data type of categorical data to factor. The attributes that will be converted to factors are Location, WindGustDir, WindDir9am, WindDir3pm, Cloud9am, Cloud3pm, and WarmerTomorrow. Now, the data frame called “new.WAUS” is ready for model fitting. The reason why I did not modify the “WAUS” data frame directly is that the assignment specification said that “WAUS” should be the default data frame name for the whole data set, therefore I think it is better to create a new data frame to store the data set after pre-processing.

Finally, I will remove all the rows with a missing value (i.e., NA). The reason is that not all the classification models can handle NAs. (e.g., Random Forest) If I only exclude the rows with NAs when building the Random Forest model (or other models that cannot handle NAs) but not for other models, this is not fair as the other models are built using some of the data that are excluded when constructing the Random Forest. Thus, to ensure consistency, I will remove all the rows with missing values before building any models, so that the data used to construct all five models are the same. Now, “new.WAUS” consists of 554 rows of data.

Question 3)

Answer:

Please refer to the code in the Appendix to generate the training and testing data set.

Question 4)

Answer:

```
# Decision Tree
# fit decision tree model to predict WarmerTomorrow
DecisionTree.new.WAUS.fit = tree(WarmerTomorrow ~ ., data=new.WAUS.train)

# Naive Bayes
# fit naive bayes model to predict WarmerTomorrow
Naive.Bayes.new.WAUS.fit=naiveBayes(WarmerTomorrow~., data = new.WAUS.train)

# Bagging
# fit bagging model to predict WarmerTomorrow
Bagging.new.WAUS.fit = bagging(WarmerTomorrow~., data = new.WAUS.train, mfinal=10)

# Boosting
# fit boosting model to predict WarmerTomorrow
Boosting.new.WAUS.fit = boosting(WarmerTomorrow~., data = new.WAUS.train, mfinal=5)

# Random Forest
# fit random forest model to predict WarmerTomorrow
Random.Forest.new.WAUS.fit = randomForest(WarmerTomorrow ~., data = new.WAUS.train)
```

Figure 3 shows the code I used to implement the five classification models

For the bagging and boosting model, I had set the mfinal to be either 5 or 10 and used the value of mfinal that gives higher accuracy. I did not use the default setting (i.e., mfinal=100) for this because the value of mfinal is too large.

Question 5)

Answer:

Please refer to the Appendix for the code I used to perform classification and create the confusion matrix

1) Decision Tree

- The confusion matrix for the decision tree classifier

```

      actual
predicted 0  1
0      53  26
1      45  43

```

Figure 4 shows the confusion matrix for decision tree model

- Accuracy = $(53+43) / (53+26+45+43) \approx 0.575 \approx 57.5\%$

2) Naïve Bayes

- The confusion matrix for the naive bayes classifier

```

      predicted
predicted 0  1
0      57  37
1      22  51

```

Figure 5 shows the confusion matrix for naive bayes classifier

- Accuracy = $(57+51) / (57+37+22+51) \approx 0.647 \approx 64.7\%$

3) Bagging

- The confusion matrix for the bagging classifier

```

      observed class
predicted class 0  1
0      39  40
1      20  68

```

Figure 6 shows the confusion matrix for bagging model

- Accuracy = $(39+68) / (39+40+20+68) \approx 0.641 \approx 64.1\%$

4) Boosting

- The confusion matrix for the Boosting classifier

```

      observed class
predicted class 0  1
0      44  30
1      35  58

```

Figure 7 shows the confusion matrix for boosting model

- Accuracy = $(44+58) / (44+30+35+58) \approx 0.611 \approx 61.1\%$

5) Random Forest

- The confusion matrix for the Random Forest classifier

```

      actual
predicted 0  1
0      40  25
1      39  63

```

Figure 8 shows the confusion matrix for random forest model

- Accuracy = $(40+63) / (40+25+39+63) \approx 0.617 \approx 61.7\%$

Question 6)

Answer:

To construct the ROC curve, I first need to calculate the confidence in predicting the target attribute for each case. After that, I will construct the ROC curve such that the curve is now showing the True Positive Rate and False Positive Rate at all confidence levels. Eventually, we will calculate the AUC value for each classifier.

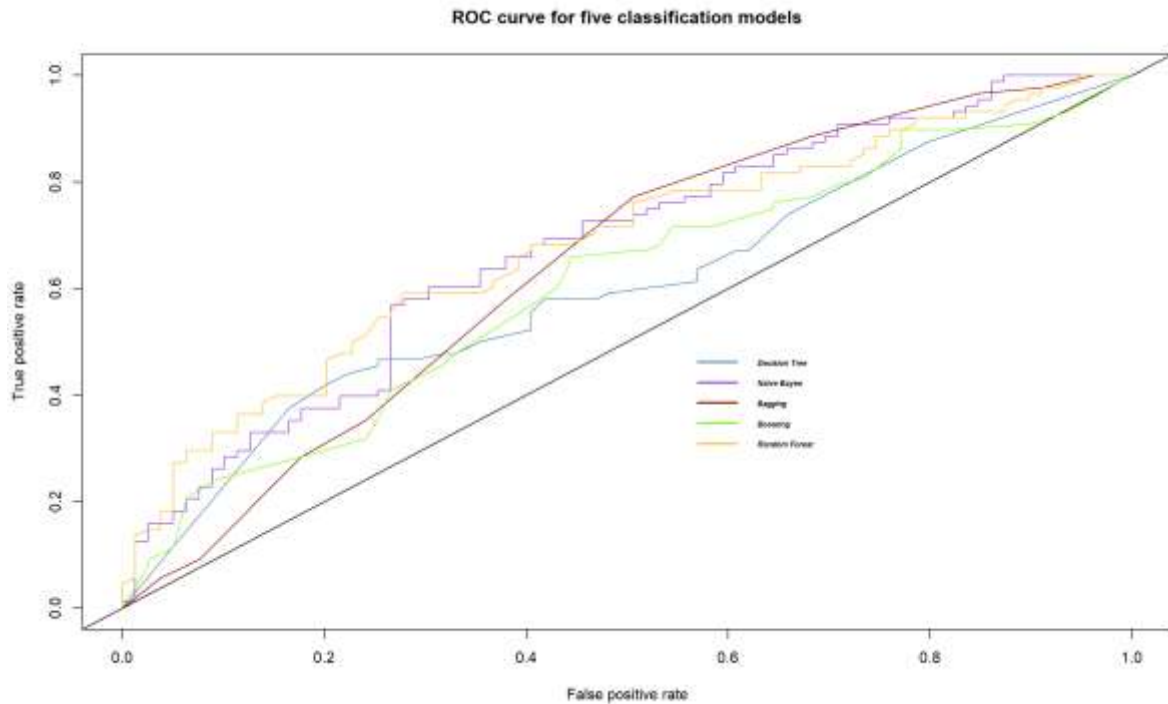


Figure 9 shows the ROC curve for each classifier

The AUC for decision tree: 0.602

The AUC for naïve bayes: 0.676

The AUC for Bagging: 0.642

The AUC for Boosting: 0.605

The AUC for Random Forest: 0.68

Question 7)

Answer:

	Decision Tree	Naïve Bayes	Bagging	Boosting	Random Forest
Accuracy	0.575	0.647	0.641	0.611	0.617
AUC	0.602	0.676	0.642	0.605	0.68

The table above shows the accuracy and AUC for all the classification models. We can see that the accuracy for the Naïve Bayes model is the highest among the 5 models, which is 0.647, followed by Bagging, Random Forest, Boosting, and Decision Tree. In addition, the

AUC value for Random Forest is the highest among the 5 models, which is 0.68, followed by Naïve Bayes, Bagging, Boosting, and Decision Tree. Higher the AUC value, the greater the ability of a classifier to differentiate between positive and negative examples over all confidence levels. In this case, the AUC values for all classifier models are between No discrimination (i.e., $AUC = 0.5$) and Acceptable discrimination (i.e., $0.7 \leq AUC < 0.8$).

By comparing these values, I would say that there is a single “best” classifier and it is the Naïve Bayes model. The reason is that it has the highest accuracy and second-highest AUC out of these models. Moreover, its AUC is just slightly lower than Random Forest at about 0.004. Thus, since the Naïve Bayes model has a better result than other models overall, I would say that the Naïve Bayes model is the “best” classifier in this situation.

Question 8)

Answer:

To answer this question, I will first determine the important or most important variables for each model, then only determine the most important variables overall. The same goes for the variables with very little effect on performance. Since there is no model for the Naïve Bayes classifier, hence we cannot determine its variable importance of it. Therefore, we will only look at the variable importance for the rest of the models.

```
#Decision Tree Attribute Importance
> print(summary(DecisionTree.new.WAUS.fit))

Classification tree:
tree(formula = warmerTomorrow ~ ., data = new.WAUS.train)
Variables actually used in tree construction:
[1] "windDir9am" "windGustDir" "Sunshine" "Humidity9am" "windDir3pm"
[6] "Humidity3pm" "Location" "Cloud3pm" "MinTemp" "Cloud9am"
[11] "Temp9am" "MaxTemp" "Evaporation" "WindSpeed9am" "Rainfall"
[16] "Pressure9am"
Number of terminal nodes: 35
Residual mean deviance: 0.456 = 161 / 352
Misclassification error rate: 0.093 = 36 / 387
```

Figure 10 shows the variables used in decision tree construction

For the decision tree model, we can determine the important variables by looking at the “variables used in tree construction” in figure 10. In other words, all the variables listed in that section can be considered the important variables to predict the target variable using the decision tree model. Thus, we know that WindGustSpeed, WindSpeed3pm, Pressure3pm, and Temp3pm do not make any contribution to the decision tree construction as they are not used in this model. In addition, the most important variables for the decision tree model are WindDir9am, WindGustDir and MaxTemp as they appear at the top branch of the decision tree model when I plot the tree. I did not include the plotting of the tree in this report as the tree is too complex to visualize. However, the code used to plot the tree is provided in the appendix.

```
#Bagging Attribute Importance
> print(Bagging.new.WAUS.fit$importance)

Cloud3pm    Cloud9am    Evaporation    Humidity3pm    Humidity9am
3.624      3.674      3.383      3.938      1.036
Location    MaxTemp    MinTemp    Pressure3pm    Pressure9am
0.850      7.229      0.749      2.972      3.013
Rainfall    Sunshine    Temp3pm    Temp9am    windDir3pm
1.974      6.179      2.051      2.794      17.187
windDir9am  windGustDir  windGustSpeed  windSpeed3pm  windSpeed9am
19.138     15.701     2.328     1.265     0.915
```

Figure 11 shows the importance of variables in predicting WarmerTomorrow for Bagging model

Figure 11 shows the importance of each attribute in the Bagging model. It shows that WindDir9am is the most important variable in the Bagging model as it has the highest value, followed by WindDir3pm and WindGustDir. Besides that, I realized that Humidity9am, Location, MinTemp, Pressure3pm, Rainfall, Temp3pm, Temp9am, WindGustSpeed, WindSpeed3pm, and WindSpeed9am, these attributes can be considered unimportant variables in this model as their important value are lesser than 3. Hence, other variables have medium importance.

```
#Boosting Attribute Importance
> print(Boosting.new.WAUS.fit$importance)
      CCloud3pm      CCloud9am      Evaporation      Humidity3pm      Humidity9am
           7.730           6.933           4.113           2.404           3.207
      Location      MaxTemp      MinTemp      Pressure3pm      Pressure9am
           1.037           3.043           0.000           2.139           0.000
      Rainfall      Sunshine      Temp3pm      Temp9am      windDir3pm
           0.958           3.344           2.169           2.139          10.316
      windDir9am      windGustDir      windGustSpeed      windSpeed3pm      windSpeed9am
          23.322          22.815           1.989           0.830           1.513
```

Figure 12 shows the importance of variables in predicting WarmerTomorrow for Boosting model

For Boosting model, it is clear to say that the important variables would be WindDir9am, WindGustDir and WindDir3pm as their importance value are relatively higher than others as shown in figure 12. On the contrary, the least important variable would be Humidity3pm, Location, MinTemp, Pressure3pm, Pressure9am, Rainfall, Temp3pm, Temp9am, WindGustSpeed, WindSpeed3pm, and WindSpeed9am. The reason is that they either do not make any contribution to the model (i.e., important is 0) or its importance value is below 3. Whereas other variables have moderate importance.

```
#Random Forest Attribute Importance
> print(Random.Forest.new.WAUS.fit$importance)
      MeanDecreaseGini
Location              3.06
MinTemp              6.51
MaxTemp              9.47
Rainfall             3.37
Evaporation          10.11
Sunshine             8.79
WindGustDir          22.07
WindGustSpeed         5.54
WindDir9am           23.58
WindDir3pm           19.44
WindSpeed9am          4.64
WindSpeed3pm          5.46
Humidity9am           7.45
Humidity3pm           9.39
Pressure9am           8.40
Pressure3pm           7.59
CCloud9am            8.31
CCloud3pm            11.44
Temp9am              7.36
Temp3pm              9.70
```

Figure 13 shows the importance of variables in predicting WarmerTomorrow for Random Forest model

For the Random Forest model, it is obvious to say that the important variables are WindGustDir, WindDir9am and WindDir3pm since their importance value are comparatively greater than others as shown in figure 13. On the other hand, I would consider Location, Rainfall, WindGustSpeed, WindSpeed9am and WindSpeed3pm as the unimportant variables to construct the Random Forest model as their importance value are lower than 6.

Overall, in conclusion, the most important variables in overall to predict WarmerTomorrow would be WindDir9am, WindDir3pm, and WindGustDir. The reason is that they are considered important variables in all/most of the classifier models. Contrastingly, the variables that could be omitted from the data with very little effect on performance in overall are WindGustSpeed and WindSpeed3pm as they are considered unimportant variables in all the four classifier models. Besides that, I would also decide to omit Pressure3pm, Temp3pm, Location, Rainfall, and WindSpeed9am as they are classified as unimportant variables in most of the classification models. In short, there are 3 important variables and 7 variables that could be omitted from the data with very little effect on performance.

Question 9)

Answer:

For this question, I decided to use create a simple tree using the decision tree model. I decided not to create any ensemble classifier (i.e., Bagging, Boosting, and Random Forest), it is because the model produced is not as easy to interpret as a single tree as many classifiers are involved during the construction of the model, hence they are not simple enough for a person to classify whether a particular data will be warmer tomorrow or not by hand. I did not consider the Naïve Bayes model because the Naïve Bayes model produced in Part 4 can be considered good enough. Unlike, the decision tree, its accuracy, and AUC are the lowest among the classifiers in Part 4. Hence, I am interested in whether making the decision tree simpler would increase its performance or not. Besides that, I also can use cross-validation to prune the decision tree down if the resulting model is complex. Thus, in my opinion, the advantages of using a decision tree outweigh the other classifiers. Above are the factors that I considered in my decision.

To build this new decision tree model, I would only use WindDir9am, WindDir3pm, and WindGustDir to predict WarmerTomorrow, the reason is that they are the most important variables determined in Part 8. Hence, I will do a similar pre-processing in Part 2 and then only create the model.

```
> table(predicted = Q9.D.predict ,actual = Q9.WAUS.test$warmerTomorrow)
      actual
predicted 0  1
      0 112 110
      1 101 142
> summary(Decision.Tree.Q9.WAUS.fit)

Classification tree:
tree(formula = WarmerTomorrow ~ ., data = Q9.WAUS.train)
Number of terminal nodes: 4
Residual mean deviance:  1.3 = 1400 / 1080
Misclassification error rate: 0.366 = 397 / 1084
> (112+142)/(112+110+101+142)
[1] 0.546
```

Figure 14 shows the accuracy, confusion matrix and summary of the new decision tree model after pre-processing the data

From the above figure, we can see that this decision tree model has 4 terminal nodes, and its accuracy is about 54.6%. Hence, I will try to prune the tree down by using a cross-validation test based on the misclassification rate to see if I can obtain a simpler tree with similar accuracy.

```
> #cross validation test
> cvtest = cv.tree(Decision.Tree.Q9.WAUS.fit, FUN = prune.misclass)
> cvtest
$size
[1] 4 3 2 1

$dev
[1] 444 444 444 499

$sk
[1] -Inf    0   28   66

$method
[1] "misclass"

attr(,"class")
[1] "prune"          "tree.sequence"
```

Figure 15 shows the details of the cross-validation test

We can see that if we prune the decision tree size to 3 (i.e., 3 terminal nodes), the dev value remains the same as the original one, and its cost complexity is 0. Thus, we can consider pruning the decision tree size to 3 to get a simpler tree.


```

> #prune using size 3 considering lowest misclassification rate
> # and lowest cost complexity
> pruned.Dfit = prune.misclass(Decision.Tree.Q9.WAUS.fit, best = 3)
> summary(pruned.Dfit)

Classification tree:
snip.tree(tree = Decision.Tree.Q9.WAUS.fit, nodes = 3L)
Variables actually used in tree construction:
[1] "WindDir9am" "WindGustDir"
Number of terminal nodes: 3
Residual mean deviance: 1.31 = 1420 / 1080
Misclassification error rate: 0.366 = 397 / 1084
> # check accuracy using the pruned tree
> PD.predict = predict(pruned.Dfit, Q9.WAUS.test, type = "class")
> table(predicted = PD.predict, actual = Q9.WAUS.test$warmerTomorrow)
      actual
predicted 0  1
      0 112 110
      1 101 142
> (112+142)/(112+110+101+142)
[1] 0.546

```

Figure 16 shows the details of the pruned decision tree model

After pruning, only WindDir9am and WindGustDir are used. Moreover, there are no changes in the accuracy and misclassification error rate between the original and pruned tree. However, a simpler tree of size 3 is created instead of the previous tree with size 4.

```

> # AUC for Question 9 Pruned Decision Tree model
> Q9.decision.tree.auc = performance(Q9.decision.tree.pred, "auc")
> print(as.numeric(Q9.decision.tree.auc@y.values)) # 0.547
[1] 0.547

```

Figure 17 shows the AUC value for pruned decision tree model

Now, we know that our final decision model has an accuracy of 54.6% and its AUC is approximately 0.547. If we compare the performance of this model with those in Part 4, I will say that this model performed the worst among all classifiers as it has the lowest accuracy and AUC.

To describe my (pruned) decision tree model, I will plot the tree and provide my explanation.



Figure 18 shows the pruned decision tree

Explanation of tree:

If the direction of the wind at 9 am (WindDir9am) is either ESE, NW, S, SE, SSE, SSW, SW, W, or WSW and the direction of the strongest wind gust over the day (WindGustDir) is either ENE, ESE, N, NNW, S, SE, SSE, SSW, SW, W, WNW, or WSW, then the model predicts that tomorrow will not be warmer than today. (WarmerTomorrow = 0)

If the direction of the wind at 9 am (WindDir9am) is either ESE, NW, S, SE, SSE, SSW, SW, W, or WSW and the direction of the strongest wind gust over the day (WindGustDir) is either NNE, NE, NW, or E, then the model predicts that tomorrow will be warmer than today. (WarmerTomorrow = 1)

If the direction of the wind at 9 am (WindDir9am) is either WNW, NNE, N, ENE, NE, NNW, or E, then the model predicts that tomorrow will be warmer than today. (WarmerTomorrow = 1)

Question 10)

For this question, I created a Bagging model and adjusted the parameters such that it gives the possible best result I can get. I will first explain why I chose the attributes I used, why I chose to create a Bagging model, what factors were important in my decision, and how I created this model, then only show the performance of my improved model.

Why I chose the attributes I used:

From Part 8, I determined what variables could be omitted from the data and the most important variables. From here, I know what variables are a must to be included in the model and what variables could be ignored. After that, I also decided to include several variables such that their importance is at least moderate in most of the classifiers as stated in Part 8. Eventually, after careful consideration by using the analysis in Part 8, I will use MaxTemp, Sunshine, WindGustDir, WindDir9am, WindDir3pm, Humidity9am, Cloud9am, Cloud3pm, Temp9am to predict WarmerTomorrow in my improved model. After that, I will do all the necessary pre-processing procedures that are similar to what I had done in Part 2 (i.e., extract selected attributes from WAUS, change the data type of categorical attribute to factor, remove rows with missing values and get training and testing set) to make sure my data set is suitable for model fitting.

Why I chose to create a Bagging model and what factors were important in my decision:

In Parts 4 and 8, I had created a total of 2 decision tree models, however, the performance is not satisfactory due to low accuracy and low AUC. Therefore, I will not consider the decision tree model anymore. Now, I am left with four types of classifiers, which are Naïve Bayes, Bagging, Boosting, and Random Forest. It is quite difficult for me to decide which classifier to use now, hence I need to fit the data to the model and check the performance. After several rounds of model fitting (i.e., adjusting the parameters, performing cross-validation, and so on), I realized that when I set mfinal to 29 for the Bagging model, I finally get the most significant improvement compared to other classifiers. Therefore, I decided to use the Bagging model as my best model for this question. I know this procedure is a bit of trial and error, but I think this is the best way to determine which models give the best result as I can try my best to not miss out on any possibility of getting a better model. Please note that I had deleted the code to construct other models and only left with the code that was used to create this Bagging model. The reason is that since I am not using those models, I think it is better to remove the piece of code to avoid confusion.

How I created my improved model:

I had set the mfinal to 29 and consider the attributes (i.e., those mentioned in the previous section) that I think can make at least a moderate contribution to predicting WarmerTomorrow for the Bagging model. Thus, eventually, the performance of the improved model (as shown in the next section) is proved to be better than the classifiers in Part 4 and 8 as its accuracy is about 0.672 (i.e., 67.2%) and AUC is around 0.698.

The performance of my best model (Bagging):

Predicted Class	Observed Class	
	0	1
0	62	36
1	26	65

Figure 19 shows the confusion matrix of the best model using Bagging

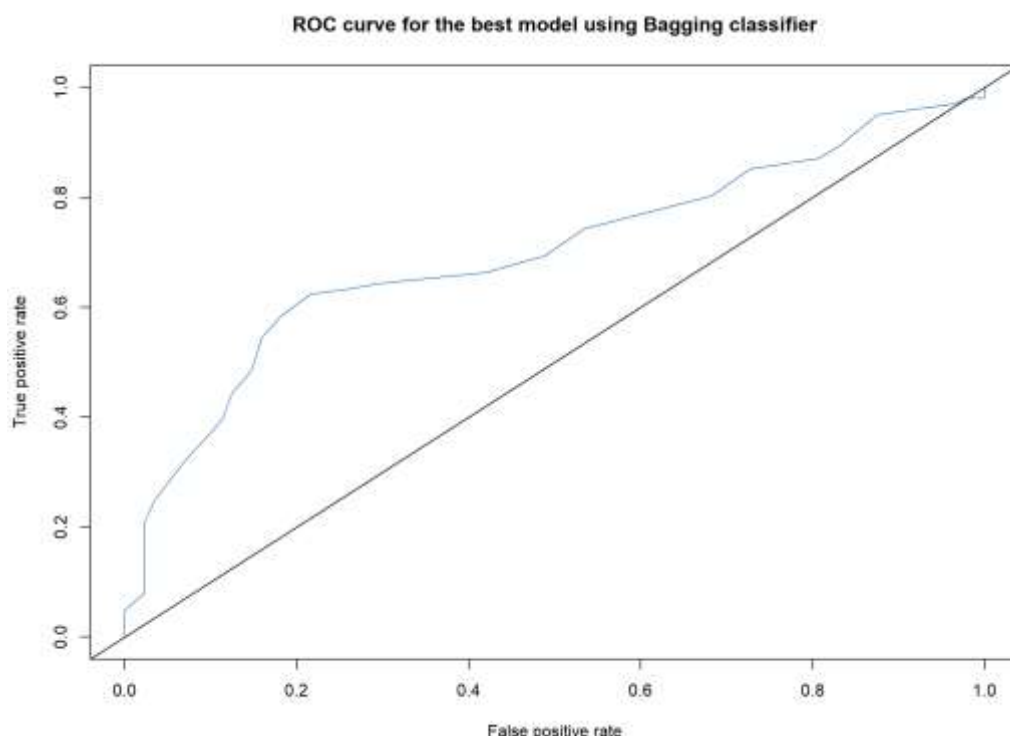


Figure 20 shows the ROC curve for the best model using Bagging classifier

```
> # AUC for Bagging model
> bagging.Q10.auc = performance(Bagging.predict.Q10.pred, "auc")
> print(as.numeric(bagging.Q10.auc@y.values)) # 0.698
[1] 0.698
```

Figure 21 shows the AUC value for the best model using Bagging classifier

From the above figures, we can know the accuracy of this model is approximately 0.672 (i.e., $(62+65) / (62+36+26+65)$) and the AUC is 0.698. Comparing this result with those classifiers in Part 4 and Part 8, we notice that the performance of this best model is better than all of them. Although the improvement is not significant (i.e., only approximately 2.5% higher for accuracy and roughly 0.02 higher for AUC than the best model in Part 4 and Part 8, which is the Naïve Bayes model in Part 4.), minor improvement can make a remarkable impact when predicting the target variable.

Question 11)

Answer:

Like Part 10, I would first comment on the attributes I used, and my data pre-processing required, then only report my Artificial Neural Network classifier performance and how it compares with the others.

Why I chose the attributes I used:

I will select MaxTemp, Sunshine, WindGustDir, WindDir9am, WindDir3pm, Humidity9am, Cloud9am, Cloud3pm and Temp9am to predict WarmerTomorrow as the attributes to construct the Artificial Neural Network classifier, which is the same as the attributes I selected in Part 10. The reason to do so is the same as I mentioned in Part 10. In short, after the analysis in Part 8, I consider these attributes could at least make a moderate contribution to the construction of the classifier. In other words, these attributes are useful to predict the value of WarmerTomorrow.

Any data pre-processing required:

The data pre-processing procedure for an Artificial Neural Network classifier is different than other models.

Here are the steps/order of the data pre-processing:

1. Extract selected attributes from "WAUS" to a new data frame called "Question.11.WAUS"
2. Remove rows containing missing values
3. Change the data type of Cloud9am and Cloud3pm to factor as they are categorical attributes so that I could use model.matrix on them
4. Create the indicator columns for each categorical attribute using model.matrix
5. Merge the output of model.matrix with "Question.11.WAUS"
6. Tidy up the merged data frame such that it only contains necessary columns
7. Separate the merged data frame into a training and testing data set
8. Fit the neural network

The performance of this Artificial Neural Network classifier:

```
> table(observed = Q11.WAUS.test$warmerTomorrow, predicted = Q11.WAUS.nn.predr)
      predicted
observed 0  1
      0 15 73
      1  5 96
> # Accuracy
> (15+96)/(15+73+5+96)
[1] 0.587
```

Figure 22 shows the confusion matrix and accuracy for Artificial Neural Network classifier

From the above figure, we know that the accuracy for this artificial Neural Network classifier is around 0.587 (i.e., 58.7%). Besides that, I had set hidden = 2 as it gives better accuracy.

Compare this classifier with the others:

This classifier has the third-lowest accuracy rate among all the classifiers in Parts 4, 9, and 10. Thus, I would not consider this classifier performing well enough to predict the attribute WarmerTomorrow.

In conclusion, the model created in Part 10 performed much better than this classifier as the accuracy of that model is around 8% higher than this Artificial Neural Network classifier.