

# Weather Forecasting in Australia

CHUN WEN KHOO

## Contents

Problem Definition.....	1
Data Preparation and Data Cleaning.....	1
Building the Models.....	2
Model Evaluation.....	2
Important Predictors.....	3
Experimentation: Parameters, Pruning, and Cross-Validation .....	3
References.....	4
Appendix.....	4

## Problem Definition

For this assignment, we have been tasked with gaining familiarity with building various classification models using R. Specifically, we would use these models to predict whether it would rain tomorrow in Australia based on ten separate locations within the country.

## Data Preparation and Data Cleaning

In the first stage of every data analysis problem, we first explore the data, in order to gain initial insight, and to discover data quality issues, if any. The initial dataset, WAUS2019.csv, has 100,000 rows. From this, I have sampled 2,000 rows of 10 random locations without replacement, using my student ID as the random seed.

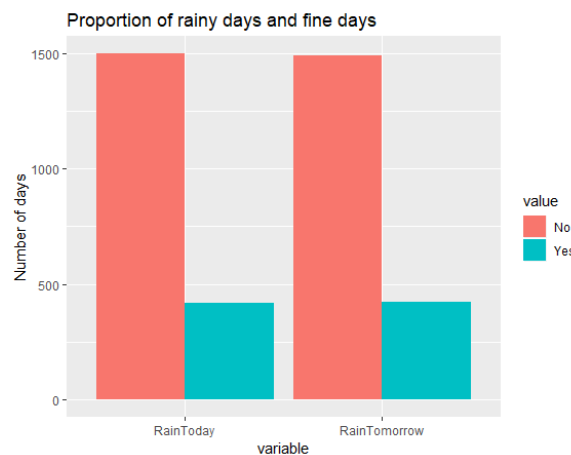


Figure 1 - Bar graph of proportion of rainy days and fine days for both actual and predicted outcomes

For starters, from the reduced dataset I will be using for the rest of this task, we can see that the proportion of actual rainy days to actual fine days is roughly 1:3.578. Thus, for every rainy day observed, we see approximately four fine days. Similarly, we can also observe that the proportion of predicted rainy days to predicted fine days is 1:3.436. This means that, for each rainy day observed, we expect there to be three, maybe four, fine days. This ties in with our earlier observation of the proportion of actual rainy days to actual fine days. It is also clear that there is an obvious class imbalance issue. However, this may just be due to the sampled locations having less rainy days.

There are two different types of data in this dataset, namely, categorical attributes, as well as numerical attributes. For the real-valued attributes, we can obtain descriptive statistics of these predictors, as seen in Figure 5. Interestingly, there are several variables with low variances, such as Evaporation and Sunshine. Thus, this may be an indication that these variables would not be good predictors. According to Alexander (2012), variables of high variance shows that there is a wide range of possibilities, both good and bad, that have a greater probability of occurring. Therefore, the converse that variables of low variance shows that there is a narrower range of possibilities, and would not be optimal for prediction, is true.

From this dataset, there are various independent variables that I believe would not be necessary to predict the weather. As an example, the variables Day, Month, and Year would not significantly affect the classification. This is because the weather does not depend on the

specific date. Thus, these variables are removed from the dataset before proceeding to the next step. Furthermore, the variable Location is also removed, due to the sole reason that we are only looking to predict if it will rain tomorrow. Hence, where it rains is irrelevant in this context. Finally, I also removed the variables Evaporation, Sunshine, Cloud9am, and Cloud3pm, since these predictors contain too many null values. As seen in Figure 4, the variables stated contain approximately 40-46% NAs, thus would not be influential in the classification. Finally, I split it into two separate sets using my student ID number as the random seed, of which 70% would be used for training the models, and the remainder would be used to evaluate the models.

## Building the Models

Using the modified dataset, I then built five different classifiers, utilizing the decision tree, the Naïve Bayes method, bagging, boosting, and random forests. All classifiers were generated using the default settings from each of their respective packages, apart from the bagging and boosting models, as these were generated using only 10 trees. Using the models built, we can then evaluate each of these classifiers based on their accuracy, confusion matrices, and the area under the receiver operating characteristic curves.

## Model Evaluation

Using the test set with each classifier, each of the test cases are classified as ‘will rain tomorrow’ or ‘will not rain tomorrow’, and these can be summarized in a confusion matrix. From their respective confusion matrices, we can then calculate the accuracy of each model, which is done by summing up the true positives and the true negatives, divided by the total number of observations. Additionally, we can also calculate the confidence of predicting ‘will rain tomorrow’ for each case and construct the ROC curve of each classifier. A ROC curve visualizes the ‘goodness’ of a classifier through graphing varying confidence threshold values of a single classifier, whereas a confusion matrix only provides comparison at a specified confidence threshold.

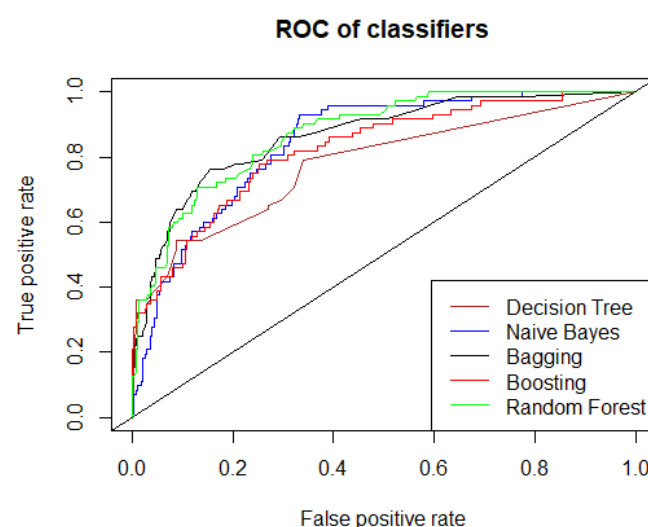


Figure 2 - Receiver Operating Characteristic curve of each classifier

From Figure 6, observe that the decision tree has the second-best accuracy, but the lowest AUC. This can be attributed to the fact that decision trees have high variance, hence one cannot be confident with the classifications of a decision tree. It can also be said that the Naïve Bayes model performs the worst, with the lowest accuracy and a relatively low AUC. This is because the Naïve Bayes method assumes predictor independence and may lose some of the underlying context. Conversely, we see that the random forest produced the highest accuracy with the highest AUC. Thus, we can conclude that the random forest best classifies this dataset.

## Important Predictors

By examining the output of each model apart from the Naïve Bayes method, we can determine the important predictors used to classify each test case. The common factors from the decision tree, bagging, boosting, and random forest methods are WindGustDir, WindDir9am, WindDir3pm, and Humidity3pm. Therefore, I believe that these variables are more important in predicting whether or not it will rain tomorrow, more so than others. Additionally, we can also remove several predictors, such as MaxTemp, RainToday, Temp9am, WindSpeed9am, and WindSpeed3pm, as these variables have been deemed unimportant by the four classifiers mentioned.

## Experimentation: Parameters, Pruning, and Cross-Validation

For starters, I performed cross validation and pruning on the initial decision tree. Based on the output, I decided to use the tree with 5 terminal nodes, as it has the lowest standard deviation. As seen in the table in Figure 7, pruning the decision tree has resulted in a higher accuracy, with a larger AUC, as compared to the original decision tree.

Additionally, I also tried increasing the final number of trees in the bagging method. By doing so, the overall accuracy of the model has increased to 84.2%, with an area under the ROC curve of 85.77%. Thus, we can see that increasing the final number of trees has slightly improved our accuracy from 83.91%, but it has also decreased the AUC by 0.78%.

Furthermore, the bagging and boosting methods can be improved upon by applying cross-validation. In this case, I have applied 10-fold cross validation, along with the bagging method, 100 final trees, and a complexity parameter of 0.01. This has resulted in a slight increase in accuracy, from 83.91% to 86.45%. I have also implemented 10-fold cross validation in the boosting methods, along with a complexity parameter of 0.01 and a Breiman learning coefficient. We see an increase in accuracy, with it jumping from 83.05% to 83.37%.

Thus, by experimenting with different parameters, I was able to generate a better classifier from the initial ones built earlier, using cross-validation with the bagging method, producing an accuracy of 86.45%.

## References

Alexander, J. (2012). Look Out for High Variance. Retrieved from <https://codermetrics.org/2012/04/11/high-variance-lookout/>

## Appendix

	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	windGustSpeed	windSpeed9am	windSpeed3pm	Humidity9am	Humidity3pm	Pressure9am
nbr.val	1976.0000	1974.0000	1950.0000	1141.0000	1082.0000	1.734e+03	1.956e+03	1.914e+03	1.935e+03	1.903e+03	1.697e+03
nbr.null	3.0000	2.0000	1307.0000	3.0000	24.0000	0.000e+00	7.300e+01	9.000e+00	0.000e+00	0.000e+00	0.000e+00
nbr.na	24.0000	26.0000	50.0000	859.0000	918.0000	2.660e+02	4.400e+01	8.600e+01	6.500e+01	9.700e+01	3.030e+02
min	-6.7000	-2.2000	0.0000	0.0000	0.0000	7.000e+00	0.000e+00	0.000e+00	2.000e+00	3.000e+00	9.923e+02
max	29.0000	46.4000	183.4000	50.2000	13.9000	9.800e+01	5.600e+01	6.500e+01	1.000e+02	1.000e+02	1.041e+03
range	35.7000	48.6000	183.4000	50.2000	13.9000	9.100e+01	5.600e+01	6.500e+01	9.800e+01	9.700e+01	4.860e+01
sum	25219.4000	45010.3000	5024.5000	6703.9000	8656.4000	7.049e+04	2.793e+04	3.595e+04	1.300e+05	1.028e+05	1.726e+06
median	12.9000	22.2000	0.0000	5.2000	9.0000	3.900e+01	1.300e+01	1.900e+01	6.900e+01	5.500e+01	1.017e+03
mean	12.7629	22.8016	2.5767	5.8755	8.0004	4.065e+01	1.428e+01	1.878e+01	6.716e+01	5.400e+01	1.017e+03
SE.mean	0.1630	0.1759	0.2244	0.1315	0.1120	3.354e-01	2.043e-01	2.027e-01	4.531e-01	4.995e-01	1.708e-01
CI.mean.0.95	0.3197	0.3450	0.4402	0.2579	0.2198	6.579e-01	4.007e-01	3.974e-01	8.887e-01	9.796e-01	3.349e-01
var	52.5240	61.0993	98.2279	19.7193	13.5744	1.951e+02	8.165e+01	7.860e+01	3.973e+02	4.748e+02	4.949e+01
std.dev	7.2473	7.8166	9.9110	4.4406	3.6844	1.397e+01	9.036e+00	8.866e+00	1.993e+01	2.179e+01	7.035e+00
coef.var	0.5678	0.3428	3.8464	0.7558	0.4605	3.436e-01	6.328e-01	4.720e-01	2.968e-01	4.035e-01	6.917e-03
nbr.val	1.696e+03	1.234e+03	1.193e+03	1935.0000	1914.0000						
nbr.null	0.000e+00	1.460e+02	9.100e+01	4.0000	2.0000						
nbr.na	3.040e+02	7.660e+02	8.070e+02	65.0000	86.0000						
min	9.869e+02	0.000e+00	0.000e+00	-5.9000	-5.1000						
max	1.038e+03	8.000e+00	8.000e+00	35.4000	45.4000						
range	5.130e+01	8.000e+00	8.000e+00	41.3000	50.5000						
sum	1.721e+06	5.317e+03	5.211e+03	33770.5000	41071.7000						
median	1.014e+03	5.000e+00	5.000e+00	17.4000	20.9000						
mean	1.014e+03	4.309e+00	4.368e+00	17.4525	21.4586						
SE.mean	1.731e-01	8.189e-02	7.958e-02	0.1689	0.1766						
CI.mean.0.95	3.396e-01	1.607e-01	1.561e-01	0.3312	0.3463						
var	5.084e+01	8.275e+00	7.555e+00	55.1997	59.6738						
std.dev	7.130e+00	2.877e+00	2.749e+00	7.4296	7.7249						
coef.var	7.028e-03	6.676e-01	6.293e-01	0.4257	0.3600						

Figure 3 - Descriptive statistics of numerical attributes of WAUS2019.csv

classifier	acc	auc
1 Decision Tree	0.8448	0.7784
2 Naive Bayes	0.8190	0.8448
3 Bagging	0.8391	0.8655
4 Boosting	0.8305	0.8233
5 Random Forest	0.8477	0.8723

Figure 4 - Comparison of accuracy and AUC of each classifier built

	ACC	AUC
original	0.8448	0.7784
cv with pruning	0.8592	0.7905

Figure 5 - Comparison of accuracy and AUC of the original decision tree and the pruned decision tree

	Observed Class	
Predicted Class	No	Yes
No	267	46
Yes	9	26

Figure 6 - Confusion matrix of bagging classifier with  $m_{final} = 100$

	Observed Class	
Predicted Class	No	Yes
No	627	90
Yes	20	75

Figure 7 - Confusion matrix of bagging classifier with cross validation

	Observed Class	
Predicted Class	No	Yes
No	608	96
Yes	39	69

Figure 8 - Confusion matrix of boosting classifier with cross validation

observed \ predicted	predicted	
	0	1
0	254	22
1	37	35

Figure 9 - Confusion matrix of the artificial neural network

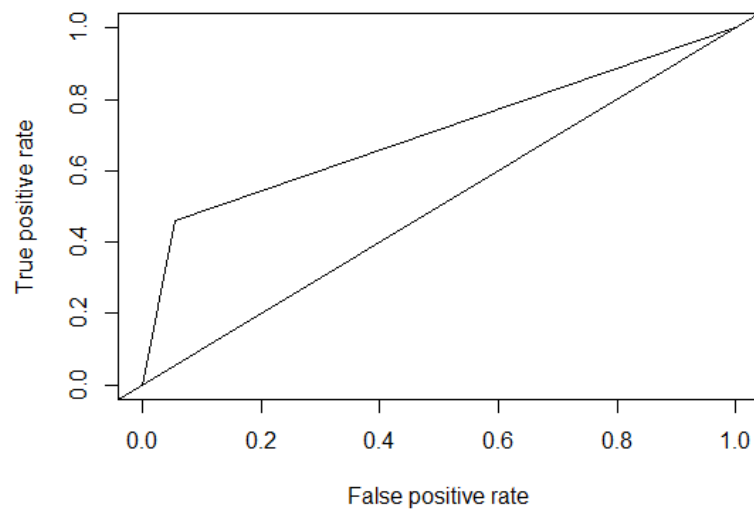


Figure 10 - ROC curve of the artificial neural network