

Rain in Australia. Classification Prediction Model

by Sumaira Afzal, Viraja Ketkar, Murlidhar Loka, Vadim Spirkov

Abstract Many native cultures comprise an institution of “rainmakers” – people who would not as much invoke the rains, but anticipate them based on ethno-meteorology. The forecasting was based on skillful art of observing the natural environment as expressed in the timing or flowering of plants, hatching of insects, arrival of migratory birds, etc., which enables farmers to make adjustments in farming calendar and crop selection types in any given season. This indigenous knowledge was often passed down from one generation to the other. We are going to employ CRISP-DM framework along with the latest scientific methods and prediction algorithms to achieve the same very goal without thorough knowledge of forces of nature, hopefully with the same accuracy as the aboriginal people.

Background

Weather forecasting is a complex and often challenging skill that involves observing and processing vast amounts of data. Weather systems can range from small, short lived thunderstorms only a few kilometers in diameter that last a couple hours to large scale rain and snow storms up to a thousand kilometers in diameter and lasting for days.

A very important component of modern weather forecasting is the use of numerical weather prediction (NWP) models. In the last years, the forecast quality of those models constantly improved, mostly due to major improvements in high performance computing. NWP focuses on taking current observations of weather and processing these data with computer models to forecast the future state of weather. Knowing the current state of the weather is just as important as the numerical computer models processing the data. Current weather observations serve as input to the numerical computer models through a process known as data assimilation to produce outputs of temperature, precipitation, and hundreds of other meteorological elements from the oceans to the top of the atmosphere.

Objective

The objective of this research is to find a supervised, binary classification model that would provide accurate forecast of the rain in Australia next day, having today's weather observations and historical data. In addition to being accurate the model should be easily interpretable and flexible enough to accept limited number of input features without diminishing its prediction power.

Data Analysis

The data set we are going to use for our research contains daily weather observations from numerous Australian weather stations collected from 2007 till 2017. There are over 142000 records. It has been sourced from [Kaggle](#)

Data Dictionary

We exclude the variable *Risk-MM* when training your binary classification model. If we don't exclude it, you will leak the answers to our model and reduce its predictability

Column Name	Column Description
Date	Date of observation
Location	Common name of the location of the weather station
MinTemp	Minimum temperature in degrees Celsius
MaxTemp	Maximum temperature in degrees Celsius
Rainfall	Amount of rainfall recorded for the day in mm
Evaporation	So-called Class A pan evaporation (mm) in the 24 hours to 9am
Sunshine	Number of hours of bright sunshine in the day

Column Name	Column Description
WindGustDir	Direction of the strongest wind gust in the 24 hours to midnight
WindGustSpeed	Speed (km/h) of the strongest wind gust in the 24 hours to midnight
WindDir9amDirection	Of the wind at 9am
WindDir3pmDirection	Of the wind at 3pm
WindSpeed9amWind	Wind speed (km/hr) averaged over 10 minutes prior to 9am
WindSpeed3pmWind	Wind Speed (km/hr) averaged over 10 minutes prior to 3pm
Humidity9amHumidity	Humidity (percent) at 9am
Humidity3pmHumidity	Humidity (percent) at 3pm
Pressure9amAtmospheric	Pressure (hpa) reduced to mean sea level at 9am
Pressure3pmAtmospheric	Pressure (hpa) reduced to mean sea level at 3pm
Cloud9amFraction	Area of sky obscured by cloud at 9am. This is measured in "oktas", which are a unit of eights. It records how many eights of the sky are obscured by cloud. A 0 measure indicates completely clear sky whilst an 8 indicates that it is completely overcast
Cloud3pmFraction	Area of sky obscured by cloud (in "oktas": eighths) at 3pm. See Cloud9am for a description of the values
Temp9amTemperature	Temperature (degrees C) at 9am
Temp3pmTemperature	Temperature (degrees C) at 3pm
RainTodayBoolean	Rainy today. 1 if precipitation (mm) in the 24 hours to 9am exceeds 1mm, otherwise 0
RISK_MM	Amount of rain. A kind of measure of the "risk". This column is redundant and will be dropped
RainTomorrow	Class label. Will it rain tomorrow?

Data Exploration

Let's take a close look at the data set. We start with loading weather observations from the file into a data frame. We remove *RISK_MM* as explained and convert *Date* column to "date" data type.

```
weatherData = read.csv("../data/weatherAUS.csv", header = TRUE, na.strings = c("NA", "", "#NA"), sep=",")
weatherData = subset(weatherData, select = -RISK_MM)
weatherData$Date = as.Date(as.character(weatherData$Date), "%Y-%m-%d")
```

Now we are going to load coordinates of the weather stations and have a bird-eye view of the weather station locations.

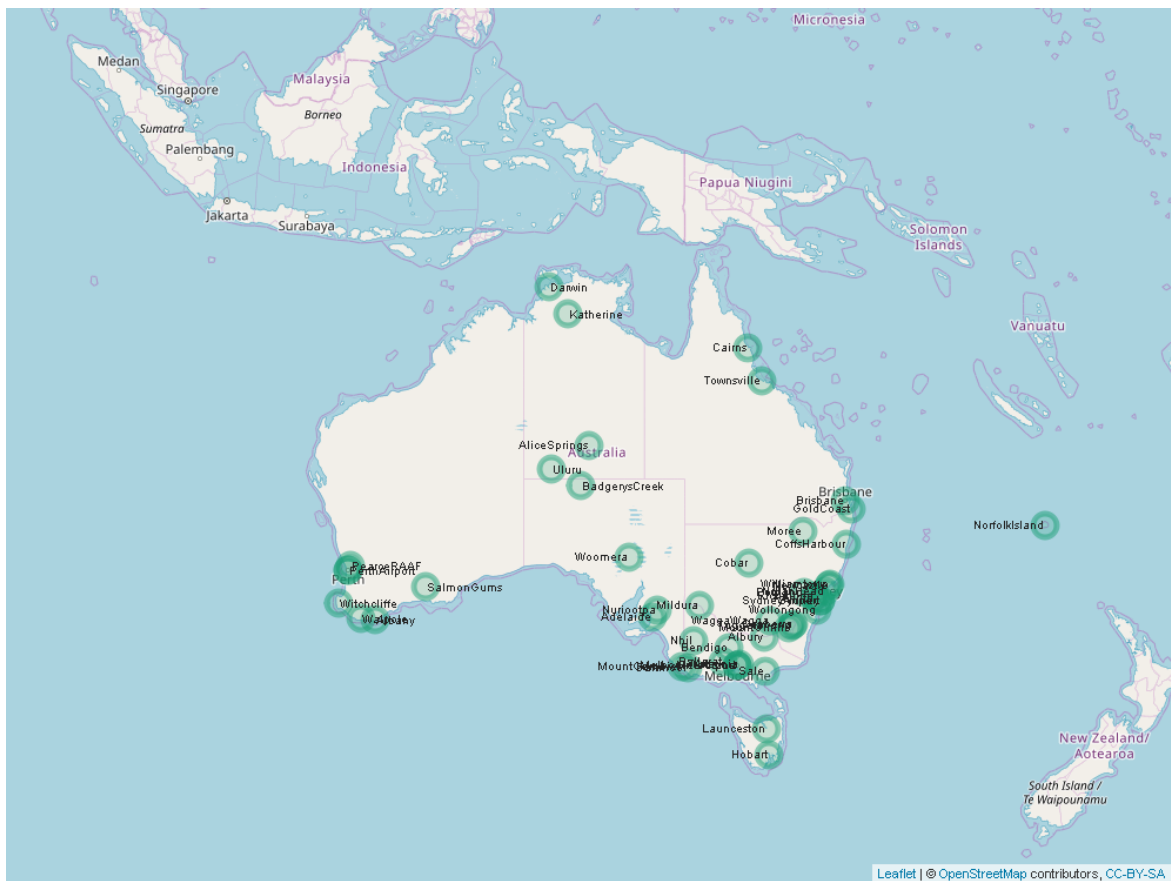


Figure 1: Australian Weather Stations

To have the full picture of the data let's print the data summary and sample.

Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir
Min. :2007-11-01	Canberra: 3418	Min. :-8.50	Min. :-4.80	Min. : 0.00	Min. : 0.00	Min. : 0.00	W : 9780
1st Qu.:2011-01-06	Sydney : 3337	1st Qu.: 7.60	1st Qu.:17.90	1st Qu.: 0.00	1st Qu.: 2.60	1st Qu.: 4.90	SE : 9309
Median :2013-05-27	Perth : 3193	Median :12.00	Median :22.60	Median : 0.00	Median : 4.80	Median : 8.50	E : 9071
Mean :2013-04-01	Darwin : 3192	Mean :12.19	Mean :23.23	Mean : 2.35	Mean : 5.47	Mean : 7.62	N : 9033
3rd Qu.:2015-06-12	Hobart : 3188	3rd Qu.:16.80	3rd Qu.:28.20	3rd Qu.: 0.80	3rd Qu.: 7.40	3rd Qu.:10.60	SSE : 8993
Max. :2017-06-25	Brisbane: 3161	Max. :33.90	Max. :48.10	Max. :371.00	Max. :145.00	Max. :14.50	(Other):86677
	(Other) :122704	NA's :637	NA's :322	NA's :1406	NA's :60843	NA's :67816	NA's : 9330

WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed9am	WindSpeed3pm	Humidity9am	Humidity3pm	Pressure9am
Min. : 6.00	N :11393	SE :10663	Min. : 0	Min. : 0.00	Min. : 0.00	Min. : 0.00	Min. : 980.5
1st Qu.: 31.00	SE : 9162	W : 9911	1st Qu.: 7	1st Qu.:13.00	1st Qu.: 57.00	1st Qu.: 37.00	1st Qu.:1012.9
Median : 39.00	E : 9024	S : 9598	Median : 13	Median :19.00	Median : 70.00	Median : 52.00	Median :1017.6
Mean : 39.98	SSE : 8966	WSW : 9329	Mean : 14	Mean :18.64	Mean : 68.84	Mean : 51.48	Mean :1017.7
3rd Qu.: 48.00	NW : 8552	SW : 9182	3rd Qu.: 19	3rd Qu.:24.00	3rd Qu.: 83.00	3rd Qu.: 66.00	3rd Qu.:1022.4
Max. :135.00	(Other):85083	(Other):89732	Max. :130	Max. :87.00	Max. :100.00	Max. :100.00	Max. :1041.0
NA's :9270	NA's :10013	NA's : 3778	NA's :1348	NA's :2630	NA's :1774	NA's :3610	NA's :14014

Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow
Min. : 977.1	Min. :0.00	Min. :0.0	Min. :-7.20	Min. :-5.40	No :109332	No :110316
1st Qu.:1010.4	1st Qu.:1.00	1st Qu.:2.0	1st Qu.:12.30	1st Qu.:16.60	Yes : 31455	Yes : 31877
Median :1015.3	Median :5.00	Median :5.0	Median :16.70	Median :21.10	NA's : 1406	
Mean :1015.3	Mean :4.44	Mean :4.5	Mean :16.99	Mean :21.69		
3rd Qu.:1020.0	3rd Qu.:7.00	3rd Qu.:7.0	3rd Qu.:21.60	3rd Qu.:26.40		
Max. :1039.6	Max. :9.00	Max. :9.0	Max. :40.20	Max. :46.70		
NA's :13981	NA's :53657	NA's :57094	NA's :904	NA's :2726		

Table 2: Weather Observations Data Summary

	Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed9am
1	14214.00	Albury	13.40	22.90	0.60			W	44	W	WNW	20
2	14215.00	Albury	7.40	25.10	0.00			WNW	44	NNW	WSW	4
3	14216.00	Albury	12.90	25.70	0.00			WSW	46	W	WSW	19
4	14217.00	Albury	9.20	28.00	0.00			NE	24	SE	E	11
5	14218.00	Albury	17.50	32.30	1.00			W	41	ENE	NW	7
6	14219.00	Albury	14.60	29.70	0.20			WNW	56	W	W	19
7	14220.00	Albury	14.30	25.00	0.00			W	50	SW	W	20
8	14221.00	Albury	7.70	26.70	0.00			W	35	SSE	W	6
9	14222.00	Albury	9.70	31.90	0.00			NNW	80	SE	NW	7
10	14223.00	Albury	13.10	30.10	1.40			W	28	S	SSE	15

WindSpeed3pm	Humidity9am	Humidity3pm	Pressure9am	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow
24	71	22	1007.70	1007.10	8		16.90	21.80	No	No
22	44	25	1010.60	1007.80			17.20	24.30	No	No
26	38	30	1007.60	1008.70		2	21.00	23.20	No	No
9	45	16	1017.60	1012.80			18.10	26.50	No	No
20	82	33	1010.80	1006.00	7	8	17.80	29.70	No	No
24	55	23	1009.20	1005.40			20.60	28.90	No	No
24	49	19	1009.60	1008.20	1		18.10	24.60	No	No
17	48	19	1013.40	1010.10			16.30	25.50	No	No
28	42	9	1008.90	1003.60			18.30	30.20	No	Yes
11	58	27	1007.00	1005.70			20.10	28.20	Yes	No

Table 3: Weather Observations Data Sample

Next set of plots renders distribution of a few important features. Overall all the features of the data set could be split into a few groups: temperature observations, humidity, wind speed, wind direction, cloud coverage, pressure, evaporation and sunshine. We picked one parameter from each group assuming that they represent well the remaining group attributes.

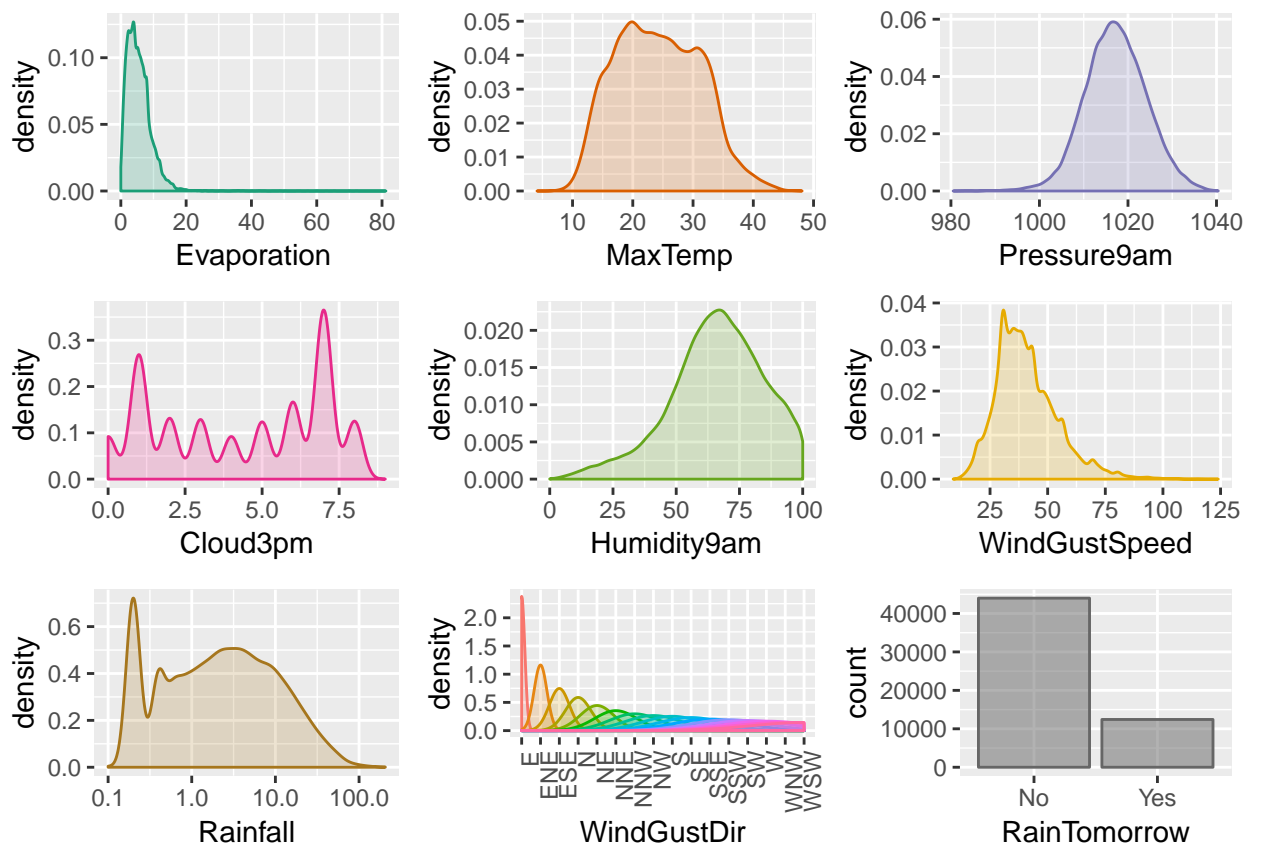


Figure 2: Distribution of Important Features

Missing Data

Data summary shows that some features miss data. Some data losses are very significant. We are going to identify what data is missing and if it is feasible to recover the missing data.

```
print(sort(colSums(is.na(weatherData)), decreasing = T), scalebox = .9)
```

Sunshine Evaporation Cloud3pm Cloud9am Pressure9am 67816 60843 57094 53657 14014 Pressure3pm WindDir9am WindGustDir WindGustSpeed WindDir3pm 13981 10013 9330 9270 3778 Humidity3pm Temp3pm WindSpeed3pm Humidity9am Rainfall 3610 2726 2630 1774 1406 RainToday WindSpeed9am Temp9am MinTemp MaxTemp 1406 1348 904 637 322 Date Location RainTomorrow 0 0 0

To speed up data processing and plot rendering we are going to use a data sample. For population of 142K observations, 30K sample size would be sufficient for 99% confidence level with the confidence interval 1.

```
weatherSample = sample_n(weatherData, SampleSize)
aggr(weatherSample, numbers = F, prop = T, col = mainPalette, sortVars = T, bars = F, varheight = T)
```

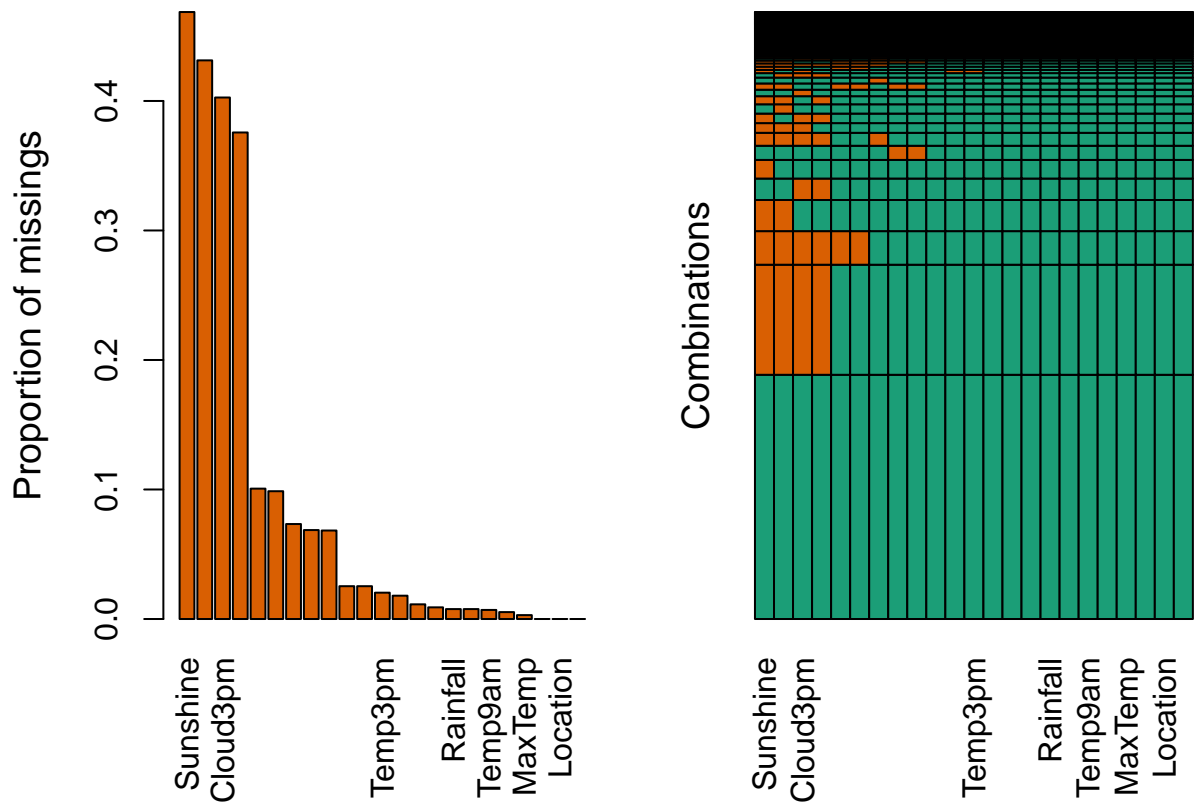


Figure 3: Missing Data Summary

```
#>
#> Variables sorted by number of missings:
#>   Variable      Count
#>   Sunshine 0.46866667
#>   Evaporation 0.43133333
#>   Cloud3pm 0.40266667
#>   Cloud9am 0.37566667
#>   Pressure9am 0.10066667
#>   Pressure3pm 0.09866667
#>   WindDir9am 0.07333333
#>   WindGustDir 0.06866667
#>   WindGustSpeed 0.06833333
#>   WindDir3pm 0.02533333
#>   Humidity3pm 0.02533333
#>   Temp3pm 0.02033333
#>   WindSpeed3pm 0.01800000
#>   Humidity9am 0.01133333
#>   WindSpeed9am 0.00900000
#>   Rainfall 0.00766667
```

```

#> RainToday 0.007666667
#> Temp9am 0.007000000
#> MinTemp 0.005333333
#> MaxTemp 0.003000000
#> Date 0.000000000
#> Location 0.000000000
#> RainTomorrow 0.000000000

```

As demonstrated in Figure 3 *Sunshine*, *Evaporation* and *Clouds* attributes suffer the loss of data between 48% and 38%. This is significant! Since we are dealing with the weather we should be observing cyclical patterns. The next group of plots illustrate missing data distribution for the most damaged features.

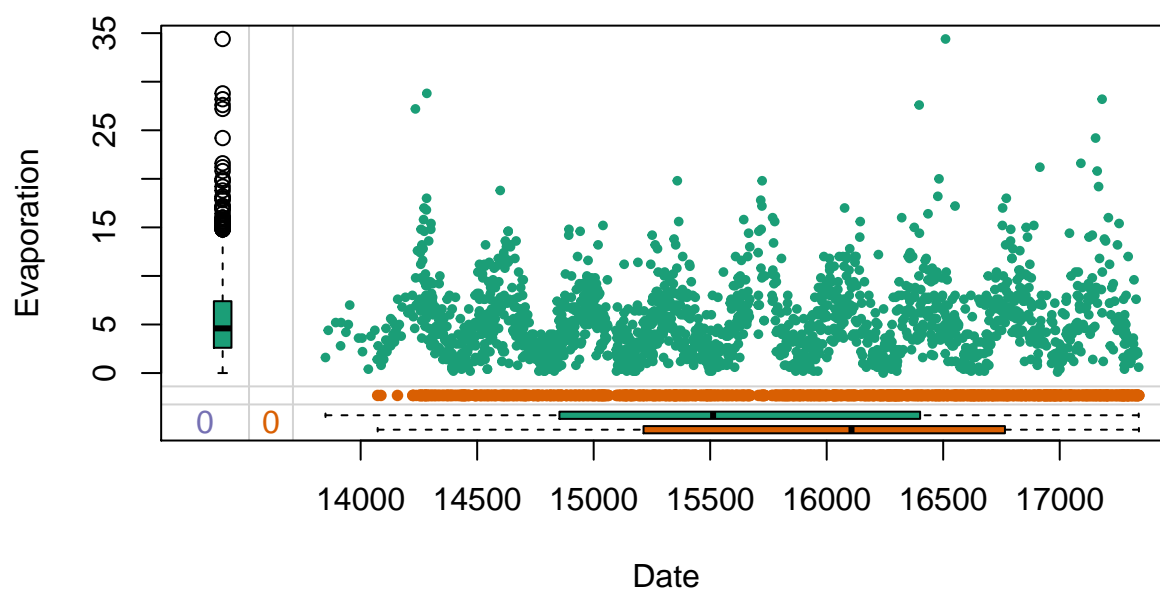


Figure 4: Date/Evaporation Margin Plot

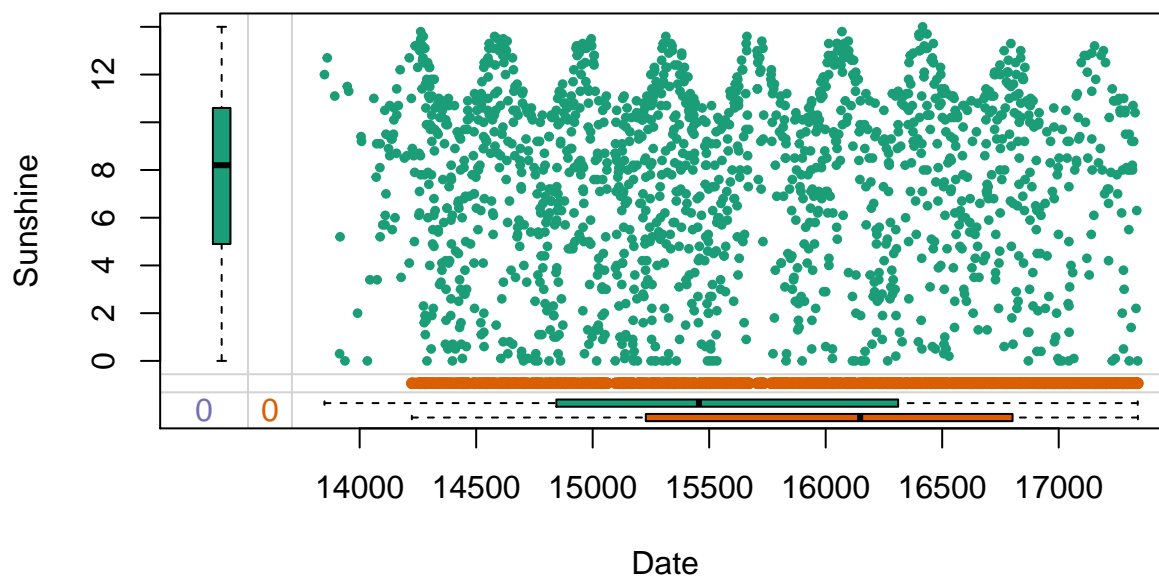


Figure 5: Date/ Sunshine Margin Plot

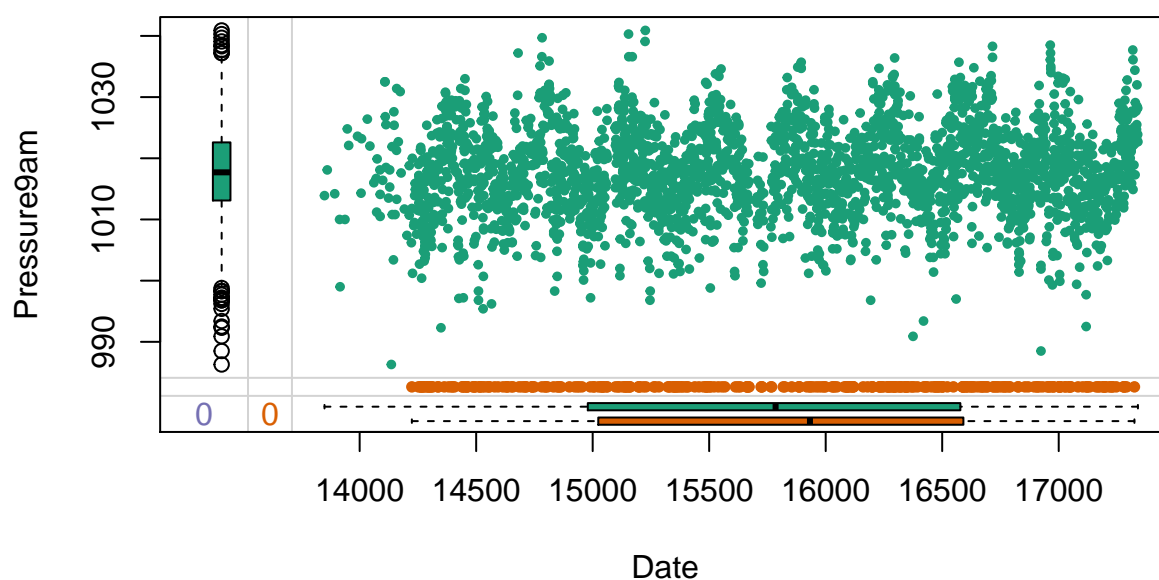


Figure 6: Date/ Pressure3pm Margin Plot

So what do the margin plots tell us? First of all let's take a look at *Date* axis. The *Date* has been converted to number to ensure continuous flow of the data. All features we picked exhibit cyclical pattern as expected. Along the vertical axis we observe the box plot of the respective feature. *Evaporation* data is quite remarkable (Figure 4); it has very narrow distribution and a lot of so-called

outliers. Though the forces of nature follow seasonal patterns they often exhibit wide range of seasonal anomalies, which the plots highlight. Thus we opt to keep the data as is. The distribution of the missing data of a given feature is depicted along the horizontal axis. In all three cases the missing data is randomly distributed over the observed date range. Along the horizontal axis we may see the box plots of the date and a given feature. Pressure readings at 9 AM *Pressure9am* (Figure 6) distributed evenly across the observed date range. *Evaporation* and *Sunshine* exhibit more data losses towards the end of the observed period.

The next plot examines another dimension of the missing data, namely missing data grouped by location.

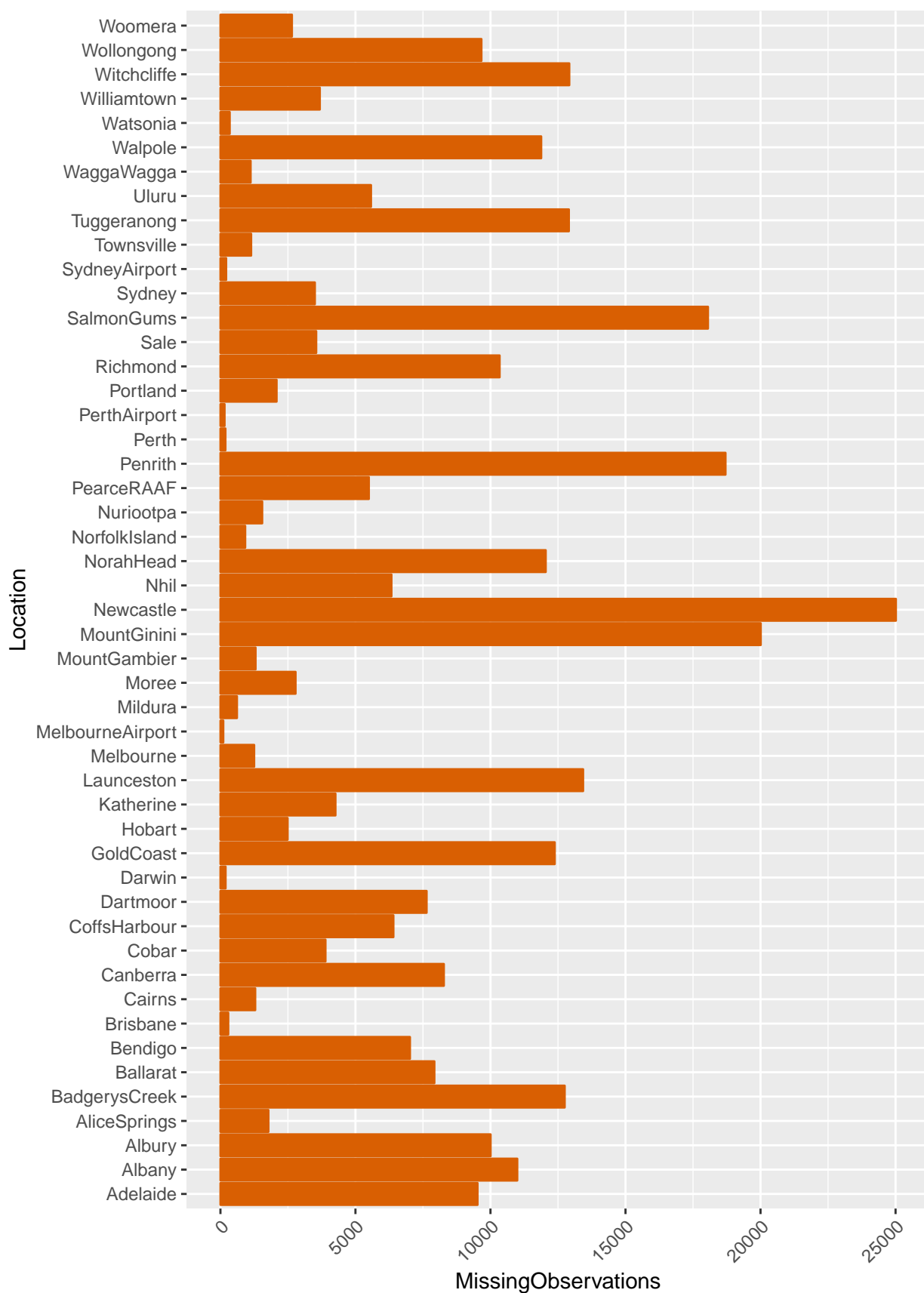


Figure 7: Missing Data By Location

Figure 7 shows that on average a location misses **6460** observations. This is significant! Though if we take a second look at the weather station map 1 we would see that Mount Gini (the station that

misses the most data), Bendigo and Ballarat are close to Melbourne, where the staff kept observing data on regular basis. Newcastle to Sydney and so on... Thus knowing that the stations are relatively close geographically we could potentially employ the weather reading collected by one station to approximate the missing data of the other station, provided they are located nearby.

Data correlation and other observations

Let's examine how the features are correlated to each other. Knowing weather we can make an assumption that the temperature features should be highly correlated, as well as pressure, wind speed, clouds and humidity groups.

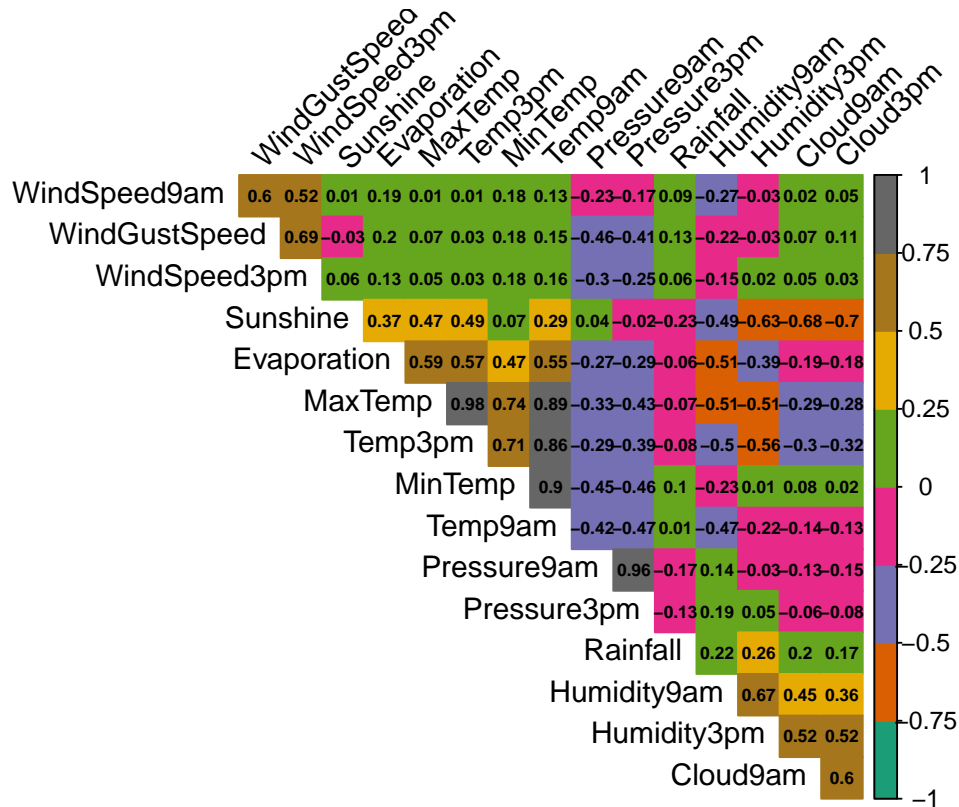


Figure 8: Data Correlation

Figure 8 confirms our initial guess. This observation will help us to eliminate redundant features later when we get to the point of selecting useful predictors for our model.

Takeaways from Data Exploration Exercise

- The data we are dealing with suffer major observation losses (Figure 3)
- The least represented features are
 - *Sunshine* **48%**.
 - *Evaporation* **43%**.
 - *Cloud* group (**40%** and **38%** respectively).
- The rest of the features exhibit medium to minor data losses, where *Pressure* group leads the way with 10%
- The missing data is distributed randomly over the observed time frame (Figures 4, 5, 6).
- We also witnessed that some weather stations recorded less data and some were almost perfect at record keeping (Figure 7). Luckily majority of the weather stations situate relatively close to each other (see figure 1). Thus if a station has data gaps the neighboring station data could be used to approximate the missing data with plausible accuracy.
- We have also noticed that many features are either positively or negatively correlated (Figure 8), where
 - *MaxTemp*, *Temp3pm* and *Temp9am* exhibit correlation of **0.86** to **0.98**.

- *Pressure9am* and *Pressure3pm* have correlation coefficient of **0.96**.
- *Sunshine* and *Cloud* group correlated negatively with coefficient of **-0.7**.
- *Rainfall* feature is of particular interest since this is what we are trying to predict. Unfortunately it does not demonstrate any strong correlations with any other features.
- Examining the data we have also seen seasonal patterns and the data that fall outside of the normal distribution range by far (outliers). Those are anomalies of nature which we opt to keep.
- The last but not least the target feature *RainTomorrow* (the value we are trying to predict) is unbalanced. So we are dealing with unbalanced data set. See Figure 2.

Data Preparation

Data exploration confirmed that despite significant data loss we should be able to impute missing data with high degree of accuracy.

Datas Imputing

To impute the missing data we employ **MICE** package. Our imputation strategy is to employ **Predictive mean matching** model. It is a robust, fast imputation algorithm that works with numeric values. Prior to applying the algorithm we have to do some data transformations, which will * increase imputation processing speed. * improve model training performance and hopefully accuracy.

First of all let's get rid of *Date* column. Outside of the presentation it does not carry too much information. What would be useful indeed is a feature that captures seasonal fluctuations. That would be *month* and *day* combined, giving us year-round (365) days of observations.

Secondly we convert *Location* categorical feature (Factor) to plain numeric column. Indeed there are 49 locations. If we dummy-encode locations how much would it improve the performance of the imputation algorithm and the processing speed (we are talking about 49 new columns...)? We believe it is more harmful than helpful. Thus we go for numeric presentation. Another topic for pondering is whether we shall employ weather station coordinates... After some deliberation we can conclude that the coordinates will not add much knowledge in the context of the model training. But they will certainly break the categorical nature of the locations. Each coordinates have 4 - 6 decimal places, which effectively makes them continuous. So we stick with numeric presentation of the locations.

We also convert other factor data types to numbers in order to make *MICE* pick **Predictive mean matching** model (vs *Multinomial logit* model, that exhibits poor speed and low performance dealing with the categorical features of 20 levels or more).

This is our original set.

```
#> 'data.frame': 142193 obs. of 23 variables:
#> $ Date : Date, format: "2008-12-01" "2008-12-02" ...
#> $ Location : Factor w/ 49 levels "Adelaide","Albany",...: 3 3 3 3 3 3 3 3 3 3 ...
#> $ MinTemp : num 13.4 7.4 12.9 9.2 17.5 14.6 14.3 7.7 9.7 13.1 ...
#> $ MaxTemp : num 22.9 25.1 25.7 28 32.3 29.7 25 26.7 31.9 30.1 ...
#> $ Rainfall : num 0.6 0 0 0 1 0.2 0 0 0 1.4 ...
#> $ Evaporation : num NA NA NA NA NA NA NA NA NA NA ...
#> $ Sunshine : num NA NA NA NA NA NA NA NA NA NA ...
#> $ WindGustDir : Factor w/ 16 levels "E","ENE","ESE",...: 14 15 16 5 14 15 14 14 7 14 ...
#> $ WindGustSpeed: int 44 44 46 24 41 56 50 35 80 28 ...
#> $ WindDir9am : Factor w/ 16 levels "E","ENE","ESE",...: 14 7 14 10 2 14 13 11 10 9 ...
#> $ WindDir3pm : Factor w/ 16 levels "E","ENE","ESE",...: 15 16 16 1 8 14 14 14 8 11 ...
#> $ WindSpeed9am : int 20 4 19 11 7 19 20 6 7 15 ...
#> $ WindSpeed3pm : int 24 22 26 9 20 24 24 17 28 11 ...
#> $ Humidity9am : int 71 44 38 45 82 55 49 48 42 58 ...
#> $ Humidity3pm : int 22 25 30 16 33 23 19 19 9 27 ...
#> $ Pressure9am : num 1008 1011 1008 1018 1011 ...
#> $ Pressure3pm : num 1007 1008 1009 1013 1006 ...
#> $ Cloud9am : int 8 NA NA NA 7 NA 1 NA NA NA ...
#> $ Cloud3pm : int NA NA 2 NA 8 NA NA NA NA NA ...
#> $ Temp9am : num 16.9 17.2 21 18.1 17.8 20.6 18.1 16.3 18.3 20.1 ...
#> $ Temp3pm : num 21.8 24.3 23.2 26.5 29.7 28.9 24.6 25.5 30.2 28.2 ...
#> $ RainToday : Factor w/ 2 levels "No","Yes": 1 1 1 1 1 1 1 1 1 2 ...
#> $ RainTomorrow : Factor w/ 2 levels "No","Yes": 1 1 1 1 1 1 1 1 2 1 ...
```

Transformation.

```
data = mutate(weatherData, MMDD = as.numeric( format(Date, "%m%d")), Location = unclass(Location),
              WindGustDir = unclass(WindGustDir),
              WindDir9am = unclass(WindDir9am), WindDir3pm = unclass(WindDir3pm),
              RainToday = unclass(RainToday)-1, RainTomorrow = unclass(RainTomorrow)-1)
data = subset(data, select = -Date)
```

Resulting data frame structure.

```
#> 'data.frame': 142193 obs. of 23 variables:
#> $ Location : int 3 3 3 3 3 3 3 3 3 3 ...
#> $ MinTemp : num 13.4 7.4 12.9 9.2 17.5 14.6 14.3 7.7 9.7 13.1 ...
#> $ MaxTemp : num 22.9 25.1 25.7 28 32.3 29.7 25 26.7 31.9 30.1 ...
#> $ Rainfall : num 0.6 0 0 0 1 0.2 0 0 0 1.4 ...
#> $ Evaporation : num NA NA NA NA NA NA NA NA NA NA ...
#> $ Sunshine : num NA NA NA NA NA NA NA NA NA NA ...
#> $ WindGustDir : int 14 15 16 5 14 15 14 14 7 14 ...
#> $ WindGustSpeed: int 44 44 46 24 41 56 50 35 80 28 ...
#> $ WindDir9am : int 14 7 14 10 2 14 13 11 10 9 ...
#> $ WindDir3pm : int 15 16 16 1 8 14 14 14 8 11 ...
#> $ WindSpeed9am : int 20 4 19 11 7 19 20 6 7 15 ...
#> $ WindSpeed3pm : int 24 22 26 9 20 24 24 17 28 11 ...
#> $ Humidity9am : int 71 44 38 45 82 55 49 48 42 58 ...
#> $ Humidity3pm : int 22 25 30 16 33 23 19 19 9 27 ...
#> $ Pressure9am : num 1008 1011 1008 1018 1011 ...
#> $ Pressure3pm : num 1007 1008 1009 1013 1006 ...
#> $ Cloud9am : int 8 NA NA NA 7 NA 1 NA NA NA ...
#> $ Cloud3pm : int NA NA 2 NA 8 NA NA NA NA NA ...
#> $ Temp9am : num 16.9 17.2 21 18.1 17.8 20.6 18.1 16.3 18.3 20.1 ...
#> $ Temp3pm : num 21.8 24.3 23.2 26.5 29.7 28.9 24.6 25.5 30.2 28.2 ...
#> $ RainToday : num 0 0 0 0 0 0 0 0 0 1 ...
#> $ RainTomorrow : num 0 0 0 0 0 0 0 1 0 ...
#> $ MMDD : num 1201 1202 1203 1204 1205 ...
```

Lets do a dry run first to see what predictors and methods for each target feature *MICE* software chooses. we will be employing **quickpred0** method of the *MICE* package. As before we will be working with a 30K data sample. Imputation process on the whole set takes about 3 hours and 20 minutes to complete!

```
meta = mice(data, maxit = 0, print = FALSE)
weatherSample = sample_n(data, SampleSize)
methods = meta$method
predictors = quickpred(data)
print(methods)
```

```
#>      Location      MinTemp      MaxTemp      Rainfall      Evaporation
#>      "      "      "pmm"      "pmm"      "pmm"      "pmm"
#>      Sunshine WindGustDir WindGustSpeed WindDir9am WindDir3pm
#>      "pmm"      "pmm"      "pmm"      "pmm"      "pmm"
#>      WindSpeed9am WindSpeed3pm Humidity9am Humidity3pm Pressure9am
#>      "pmm"      "pmm"      "pmm"      "pmm"      "pmm"
#>      Pressure3pm Cloud9am      Cloud3pm      Temp9am      Temp3pm
#>      "pmm"      "pmm"      "pmm"      "pmm"      "pmm"
#>      RainToday RainTomorrow      MMDD
#>      "pmm"      "      "      "      "
```

The code output above shows that: 1. the features without missing data will not be imputed. 2. The imputation targets will all be treated with *Predictive mean matching* algorithm ("pmm").

This is exactly what we need. Now let's review the predictors (*Code output is not included into report to save space*).

The matrix of predictors has the predictors in the columns and the features to be imputed in the rows. If the cell value equals **1** the predictor will be employed in calculations of the respective imputation target. Surprisingly *MMDD* is not used widely to predict the missing data, neither do the *Location*.

Now we are going to start the imputation process. **Note: it might take about 4 - 5 minutes even for a sample.** We have disabled the output of the function as we do not want to pollute the report with irrelevant messages

```
imputed = mice(weatherSample, pred = predictors, meth = methods, seed = 38019,
               nnet.MaxNWts = 2000, printFlag = F)
```

Now it is time to analyze the imputed values. In general, a good imputed value is a value that could have been observed had it not been missing. The MAR assumption can never be tested from the observed data. To check whether the imputations created by **MICE** algorithm are plausible we employ density charts and compare the distribution of the imputed values vs real observations. Let's do this (*again the plots take time, patience...*).

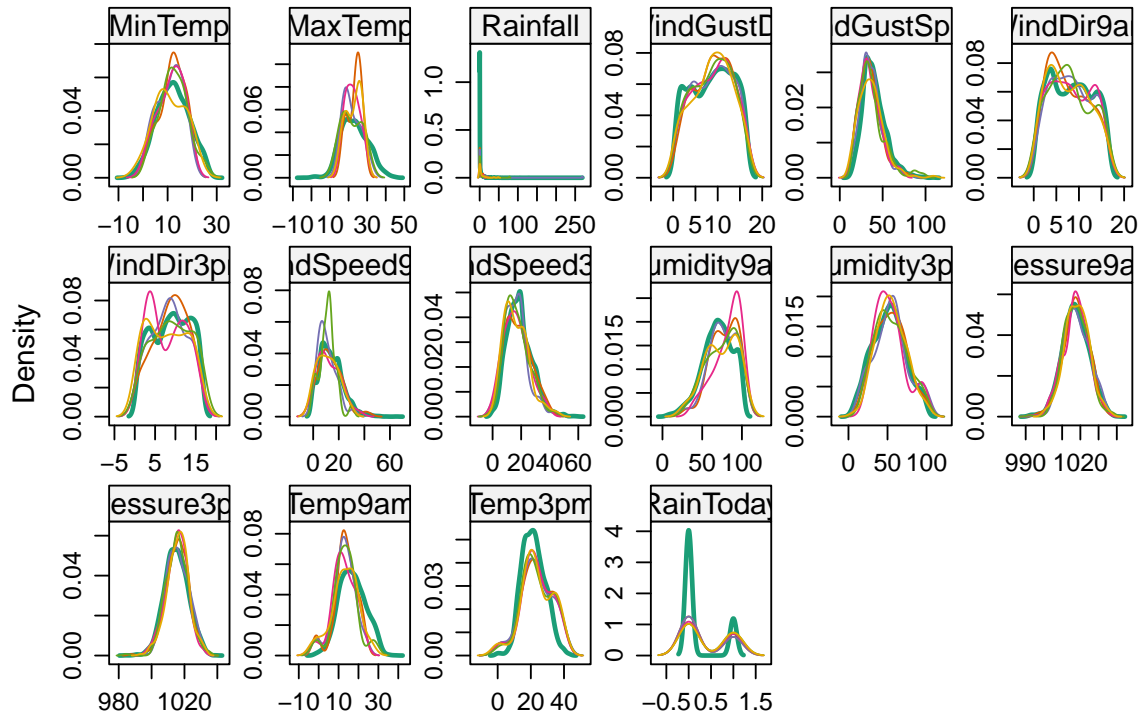


Figure 9: Imputed Values Distribution vs Real Observations

Figure 9 illustrates imputed value distribution for each imputed feature vs observed data. The fat green line renders the real data distribution and the thin lines of other colors show the distribution of the imputed data after each imputation cycle (*there are five of them by default*). The last imputation run is rendered in yellow; it should be shadowing the contour of the green one as close as possible. And it does which give us an indication that the result of the imputation is plausible. So looking at the charts we can conclude that the imputation has been successful! Let's apply imputed values to our sample set and verify if there are any *NAs* left.

```
weatherSample = complete(imputed)
print(colSums(is.na(weatherSample)))
```

```
#>      Location      MinTemp      MaxTemp      Rainfall      Evaporation
#>      0              0          0          0              0
#>      Sunshine  WindGustDir WindGustSpeed WindDir9am  WindDir3pm
#>      0              0          0          0              0
#>      WindSpeed9am WindSpeed3pm Humidity9am Humidity3pm  Pressure9am
#>      0              0          0          0              0
#>      Pressure3pm  Cloud9am    Cloud3pm    Temp9am     Temp3pm
#>      0              0          0          0              0
#>      RainToday   RainTomorrow      MMDD
#>      0              0          0
```

Outstanding! There are no missing values. Now we move on to the next part - model training.

Modeling and Evaluation

Finally we have reached the stage where we can start training and evaluating classification models. At this point we have clear understanding of our data. We have gotten rid of the features that did not present much value. We have filled the gaps in our data set employing sophisticated imputation technique.

Feature Selection

The weather observation data set originally had 24 features. We have removed *RISK_MM* and *Date* as explained earlier and added *MMDD*. Now the data set has 22 features and one label - *RainTomorrow*. Let's see if we can reduce the number of predictors without significant information loss. This would make our models faster and more interpretable for users. We shall keep in mind that at the data exploration phase we discovered that many features were correlated (Figure 8). Hopefully this knowledge will help us to identify and remove redundant features.

Generally speaking feature evaluation methods can be separated into two groups: those that use the model information and those that do not. Clearly at this stage of our research the models are not ready. Thus we will be exploring the methods that do not require model.

This group of the method could be spit further as follows: * wrapper methods that evaluate multiple models adding and/or removing predictors. These are some examples: + recursive feature elimination + genetic algorithms + simulated annealing

- filter methods which evaluate the relevance of the predictors outside of the predictive models.

The evaluation of various feature selection methods is not in the scope of this paper. Thus we opt for a recursive feature elimination method using accuracy as a target metric.

Before we precede any further let's ensure that all categorical values get converted back to the factors. This is useful for dimentionality reduction algorithms and model training.

```
weatherSample = mutate(weatherSample, Location = as.factor(unclass(Location)),
  WindGustDir = as.factor(unclass(WindGustDir)),
  WindDir9am = as.factor(unclass(WindDir9am)), WindDir3pm = as.factor(unclass(WindDir3pm)),
  RainToday = as.factor(unclass(RainToday)), RainTomorrow = as.factor(unclass(RainTomorrow)))
```

It is time to run feature selection algorithm.

```
predictors = subset(weatherSample, select = -RainTomorrow)
label = weatherSample[,22]

# run the RFE algorithm
rfePrediction = rfe(predictors, label, sizes=c(1:22),
  rfeControl = rfeControl(functions=rffuncs, method="cv", number=3))
print(rfePrediction)

#>
#> Recursive feature selection
#>
#> Outer resampling method: Cross-Validated (3 fold)
#>
#> Resampling performance over subset size:
#>
#> Variables Accuracy Kappa AccuracySD KappaSD Selected
#>      1  0.8017 0.3532 0.0102144 0.046680
#>      2  0.7813 0.3474 0.0085049 0.031708
#>      3  0.8083 0.4046 0.0102144 0.042480
#>      4  0.8000 0.3967 0.0020000 0.021957
#>      5  0.8127 0.4344 0.0075056 0.029696
#>      6  0.8137 0.4364 0.0065064 0.014101
#>      7  0.8203 0.4696 0.0070238 0.020086
#>      8  0.8223 0.4710 0.0080208 0.024674
#>      9  0.8237 0.4853 0.0030551 0.039928
#>     10  0.8210 0.4882 0.0036056 0.021525
#>     11  0.8210 0.4872 0.0055678 0.024823
#>     12  0.8203 0.4896 0.0030551 0.021626
```

```
#>      13  0.8227 0.4985 0.0025166 0.021794
#>      14  0.8270 0.5045 0.0020000 0.027149      *
#>      15  0.8247 0.5021 0.0055076 0.029916
#>      16  0.8247 0.5004 0.0005774 0.017920
#>      17  0.8253 0.5032 0.0037859 0.021247
#>      18  0.8247 0.5007 0.0005774 0.021494
#>      19  0.8267 0.5061 0.0030551 0.015951
#>      20  0.8250 0.5045 0.0043589 0.020621
#>      21  0.8257 0.5026 0.0049329 0.019944
#>      22  0.8243 0.4985 0.0051316 0.008843
#>
#> The top 5 variables (out of 14):
#>      Humidity3pm, Sunshine, Cloud3pm, WindGustSpeed, Pressure3pm
```

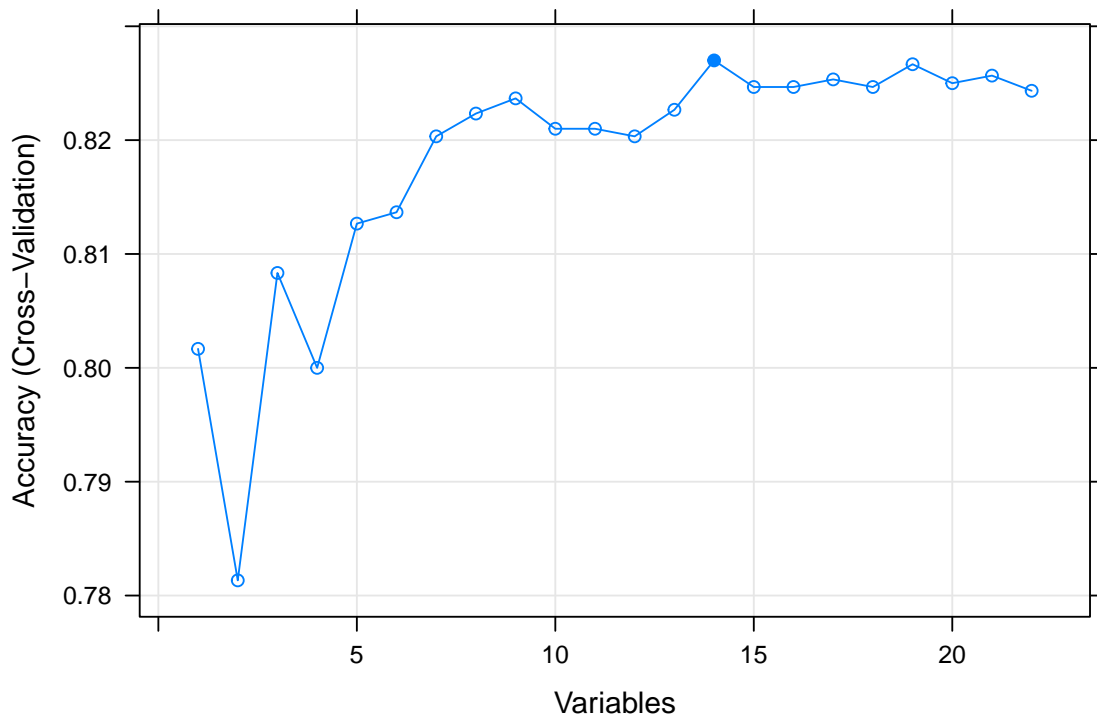


Figure 10: Number of Predictors vs Accuracy

Figure 10 suggests that the accuracy peaks a few times: with 9 predictors then 14 and tops at 22. The accuracy gain between 9 and 22 is negligible. Here is the list of features ordered by importance. We take first nine for model training.

```
#> [1] "Humidity3pm" "Sunshine" "Cloud3pm" "WindGustSpeed"
#> [5] "Pressure3pm" "Rainfall" "Pressure9am" "Location"
#> [9] "Cloud9am" "Humidity9am" "RainToday" "WindSpeed3pm"
#> [13] "MinTemp" "Temp3pm"
```

It is counterintuitive that none of the temperature readings made it to the top, neither did our synthetic *MMDD* feature. Biases destroyed!

Data Upsampling

There is one more step to make before we get to the model training. As shown in Figure 2 our data set is unbalanced. This could cause model over-fitting. So let's split the data into the training and testing sets and up-sample the training set.

```
set.seed(1608)
```

```
# keep only the selected features
finalSample = weatherSample %>% dplyr::select(c(selectedPredictors,"RainTomorrow"));

splitIdx = createDataPartition(finalSample$RainTomorrow, p=0.7, list = F) # 70% training data
trainData = finalSample[splitIdx, ]
testData = finalSample[-splitIdx, ]

set.seed(590045)
columns = colnames(trainData)
trainData = upSample(x = trainData[, columns[columns != "RainTomorrow"] ],
  y = trainData$RainTomorrow, list = F, yname = "RainTomorrow")

rm(splitIdx, columns, finalSample)
print(table(trainData$RainTomorrow))

#>
#>    0    1
#> 1597 1597
```

As we can see now the training set is balanced.

Thus we have prepared our training and test data sets. We have identified the most important features. We are ready to work on the prediction models.

Classification (Decision) Tree Model

Decision trees tend to be the method of choice for predictive modeling. A classification tree is used to predict qualitative data.

```
#> Conditional Inference Tree
#>
#> 3194 samples
#>    9 predictor
#>    2 classes: 'no', 'yes'
#>
#> No pre-processing
#> Resampling: Cross-Validated (5 fold)
#> Summary of sample sizes: 2555, 2555, 2556, 2555, 2555
#> Resampling results across tuning parameters:
#>
#>   mincriterion   ROC       Sens       Spec
#>   0.01          0.8502123 0.7620141 0.8002978
#>   0.50          0.8450125 0.7695239 0.7727253
#>   0.99          0.8263517 0.7457935 0.7540106
#>
#> ROC was used to select the optimal model using the largest value.
#> The final value used for the model was mincriterion = 0.01.

confusionMatrix(data = pred.classTreeModel.raw, testDataCopy$RainTomorrow)

#> Confusion Matrix and Statistics
#>
#>              Reference
#> Prediction no yes
#>      no    562  69
#>      yes   122 146
#>
#>              Accuracy : 0.7875
#>              95% CI : (0.7593, 0.8139)
#>      No Information Rate : 0.7608
#>      P-Value [Acc > NIR] : 0.0318239
#>
#>              Kappa : 0.4617
#>      McNemar's Test P-Value : 0.0001682
#>
#>              Sensitivity : 0.8216
```



```
#>      Specificity : 0.6791
#>      Pos Pred Value : 0.8906
#>      Neg Pred Value : 0.5448
#>      Prevalence : 0.7608
#>      Detection Rate : 0.6251
#>      Detection Prevalence : 0.7019
#>      Balanced Accuracy : 0.7504
#>
#>      'Positive' Class : no
#>
```

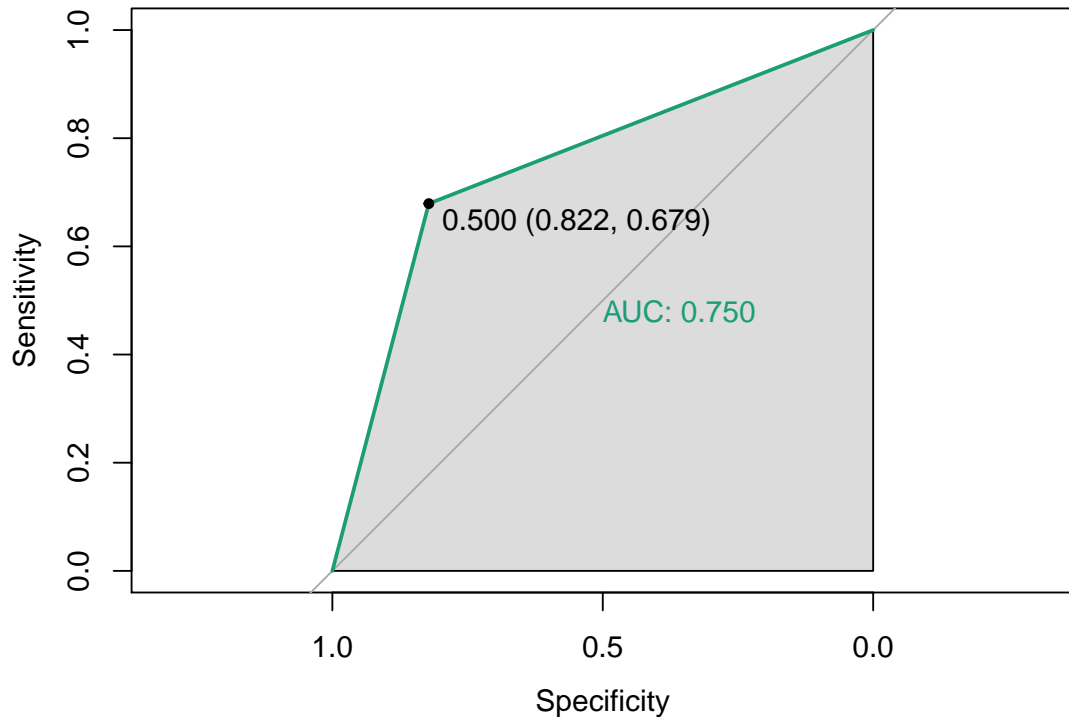


Figure 11: Classification Tree Model AUC and ROC Curve

Naive Bayes Model

```
#> Naive Bayes
#>
#> 3194 samples
#> 9 predictor
#> 2 classes: 'no', 'yes'
#>
#> No pre-processing
#> Resampling: Cross-Validated (5 fold)
#> Summary of sample sizes: 2555, 2555, 2556, 2555, 2555
#> Resampling results across tuning parameters:
#>
#> usekernel ROC Sens Spec
#> FALSE 0.5000000 NaN NaN
#> TRUE 0.8284473 0.8672218 0.6117947
#>
#> Tuning parameter 'fL' was held constant at a value of 0
#> Tuning
#> parameter 'adjust' was held constant at a value of 1
#> ROC was used to select the optimal model using the largest value.
```

```
#> The final values used for the model were fL = 0, usekernel = TRUE
#> and adjust = 1.

confusionMatrix(data = pred.naiveBayesModel.raw, testDataCopy$RainTomorrow)

#> Confusion Matrix and Statistics
#>
#>           Reference
#> Prediction no yes
#>      no  597  86
#>      yes   87 129
#>
#>              Accuracy : 0.8076
#>              95% CI : (0.7802, 0.8328)
#>      No Information Rate : 0.7608
#>      P-Value [Acc > NIR] : 0.000455
#>
#>              Kappa : 0.4721
#>  Mcnemar's Test P-Value : 1.000000
#>
#>              Sensitivity : 0.8728
#>              Specificity : 0.6000
#>      Pos Pred Value : 0.8741
#>      Neg Pred Value : 0.5972
#>      Prevalence : 0.7608
#>      Detection Rate : 0.6641
#>      Detection Prevalence : 0.7597
#>      Balanced Accuracy : 0.7364
#>
#>      'Positive' Class : no
#>
```

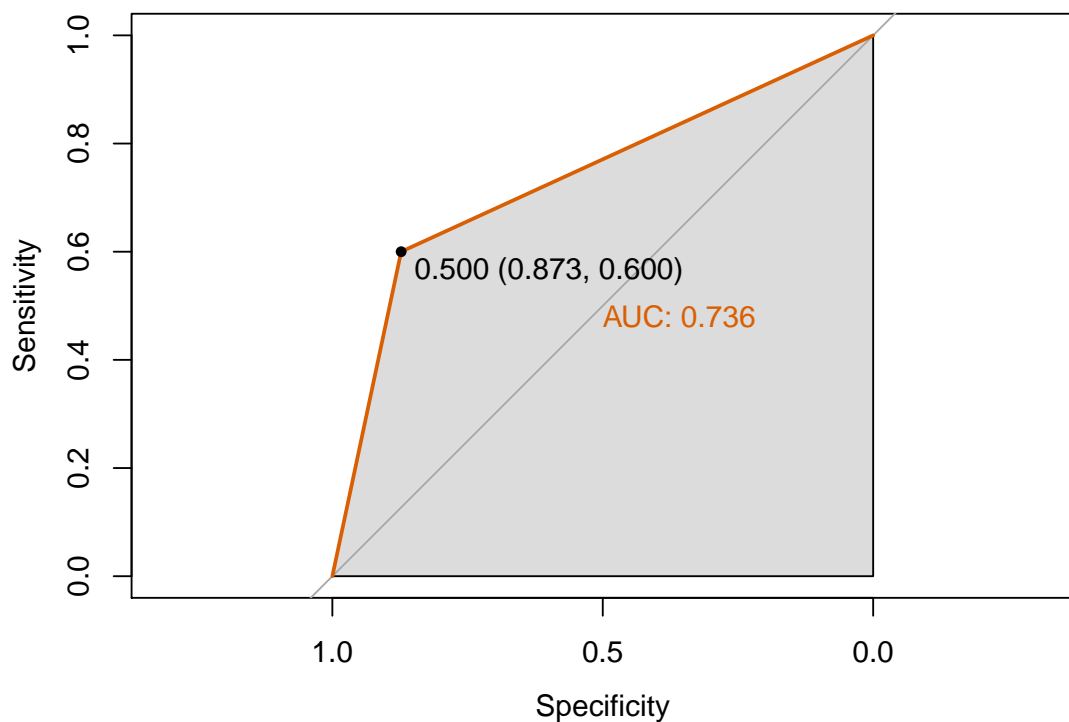


Figure 12: Naive Bayes Model AUC and ROC Curve

Random Forest Model

```

#> Random Forest
#>
#> 3194 samples
#>    9 predictor
#>    2 classes: 'no', 'yes'
#>
#> No pre-processing
#> Resampling: Cross-Validated (2 fold)
#> Summary of sample sizes: 1597, 1597
#> Resampling results across tuning parameters:
#>
#>  mtry  ROC          Sens          Spec
#>    2   0.8802078  0.7827304  0.7964882
#>   29   0.9580534  0.8422111  0.9405091
#>   56   0.9547915  0.8359486  0.9367529
#>
#> ROC was used to select the optimal model using the largest value.
#> The final value used for the model was mtry = 29.

confusionMatrix(data = pred.randomForestModel.raw, testDataCopy$RainTomorrow)

#> Confusion Matrix and Statistics
#>
#>              Reference
#> Prediction  no  yes
#>          no  615  87
#>          yes  69 128
#>
#>              Accuracy : 0.8265
#>              95% CI : (0.8001, 0.8507)
#>    No Information Rate : 0.7608
#>    P-Value [Acc > NIR] : 1.111e-06
#>
#>              Kappa : 0.5091
#>  McNemar's Test P-Value : 0.1735
#>
#>              Sensitivity : 0.8991
#>              Specificity : 0.5953
#>              Pos Pred Value : 0.8761
#>              Neg Pred Value : 0.6497
#>              Prevalence : 0.7608
#>              Detection Rate : 0.6841
#>    Detection Prevalence : 0.7809
#>              Balanced Accuracy : 0.7472
#>
#>              'Positive' Class : no
#>

```

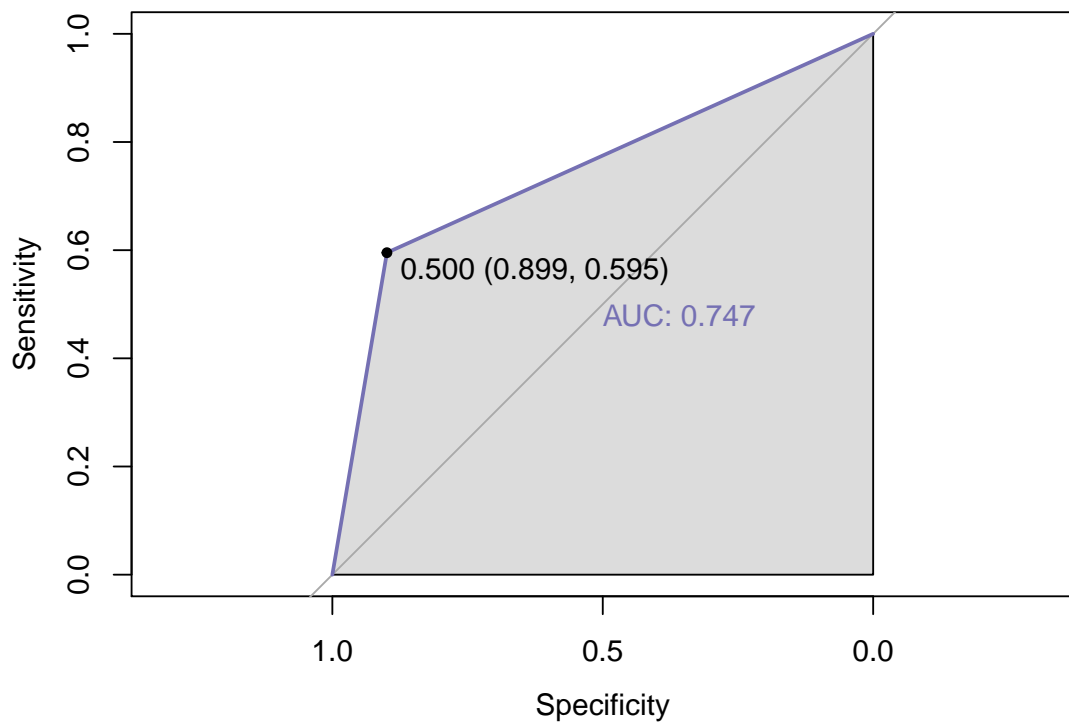


Figure 13: Random Forest Model AUC and ROC Curve

Logistic Regression Model

Logistic regression is an efficient, interpretable and accurate method, which fits quickly with minimal tuning. Logistic regression prediction accuracy will benefit if the data is close to Gaussian distribution. Thus we apply addition transformation to the training data set. We will also be employing 5-fold cross-validation resampling procedure to improve the model. In addition to the above we are going to convert *Location* categorical value to numeric data type. We could have used dummy encoding but having 49 locations such approach does not seem beneficial.

```
confusionMatrix(data = pred.logRegModel.raw, testDataCopy$RainTomorrow)
```

```
#> Confusion Matrix and Statistics
#>
#>      Reference
#> Prediction  0   1
#>      0  538  46
#>      1  146 169
#>
#>              Accuracy : 0.7864
#>              95% CI : (0.7582, 0.8128)
#>      No Information Rate : 0.7608
#>      P-Value [Acc > NIR] : 0.038
#>
#>              Kappa : 0.4938
#>  Mcnemar's Test P-Value : 9.019e-13
#>
#>      Sensitivity : 0.7865
#>      Specificity : 0.7860
#>      Pos Pred Value : 0.9212
#>      Neg Pred Value : 0.5365
#>      Prevalence : 0.7608
#>      Detection Rate : 0.5984
#>      Detection Prevalence : 0.6496
#>      Balanced Accuracy : 0.7863
```

```
#>
#> 'Positive' Class : 0
#>
```

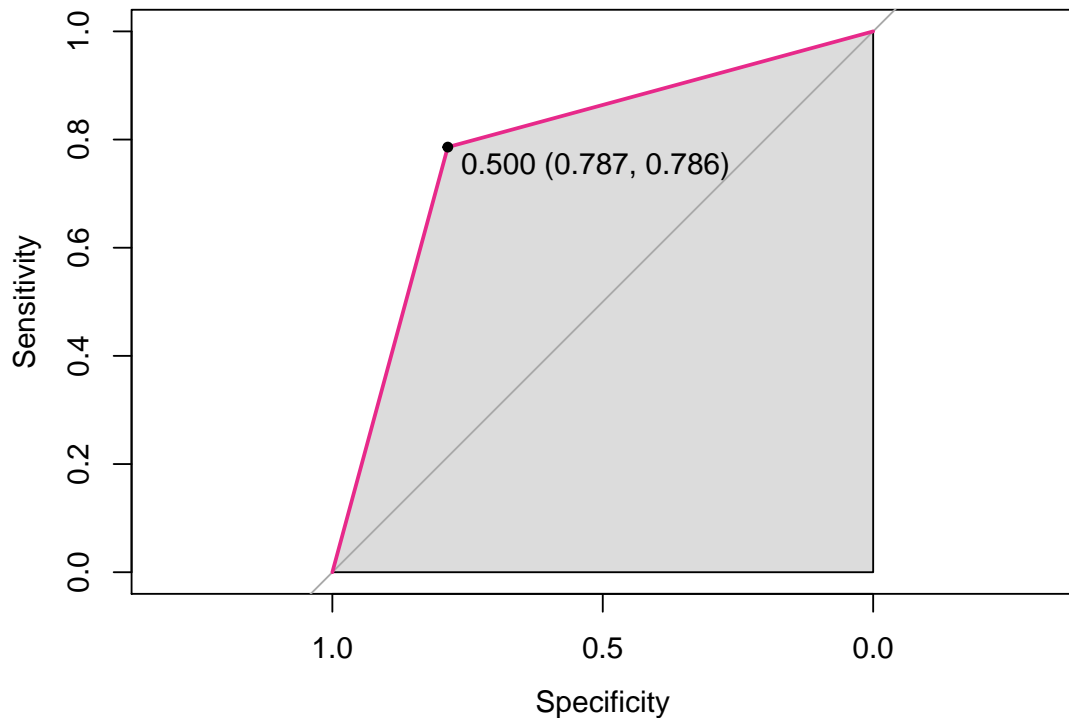


Figure 14: Logistic Regression Model AUC and ROC Curve

Confusion matrix and Figure 14 demonstrate the logistic model performance on the balanced data set. Using the proportion of positive data points that are correctly considered as positive (true positives) and the proportion of negative data points that are mistakenly considered as positive (false negative), we generated a graphic that shows the trade off between the rate at which the model correctly predicts the rain tomorrow with the rate of incorrectly predicting the rain. The value around 0.80 indicates that the model does a good job in discriminating between the two categories.

Model Comparison

Now it is time to compare the models side by side and pick a winner.

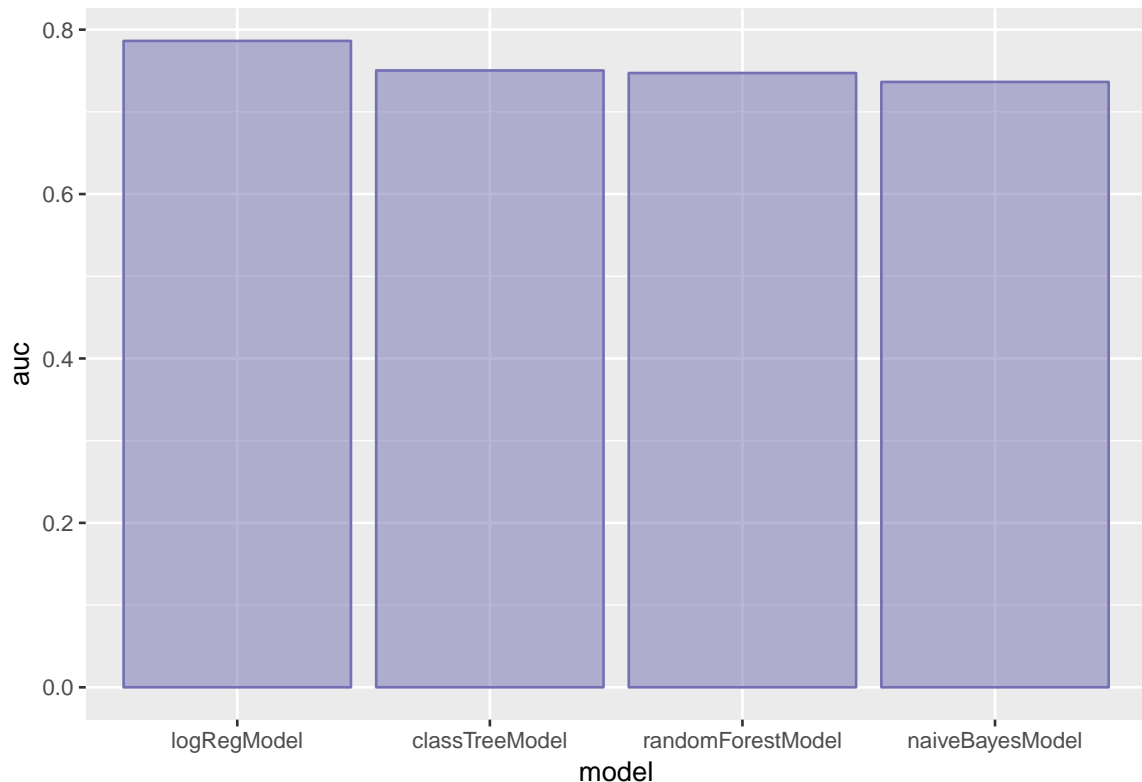


Figure 15: Model AUC Comparison

```
#>           model      auc
#> 1    logRegModel 0.7862981
#> 2   classTreeModel 0.7503536
#> 4 randomForestModel 0.7472358
#> 3   naiveBayesModel 0.7364035
```

AUC - ROC performance AUC stands for Area under the ROC Curve and ROC for Receiver operating characteristic curve. This is one of the most important KPIs of the classification algorithms. These two metrics measure how well the models distinguishing between the classes. The higher AUC the better model predicts positive and negative outcome.

Figures 11, 12, 13, 14 and accompanying data show that on the test data set the Random Forest model has the higher overall accuracy (83%) but performs poorly predicting rainy days (59%), thus the balanced accuracy is lower (about 75%).

Naive Bayes and Decision Tree are less accurate models in comparison with the Random Forest one but way more balanced demonstrating consistent power to predict rainy and sunny days with almost equal accuracy. They have balanced accuracy of 74% and 79% respectively.

Logistic regression model scores the best having the highest AUC and all other metrics around 80%.

Model interpretability Logistic Regression, Decision Tree and Naive Bayes are all highly interpretable models. It is easy to explain to the business what impact each input parameter has. The decision tree could be visualized (provided if it is not too large).

Random Forest on the other hand is a black-box model, complex algorithm which is difficult to explain in simple terms.

Data Preparation Decision Tree, Random Forest and Naive Bayes can deal with missing data, outliers, numeric and alphanumeric values. Simply speaking they are not very demanding for data quality. It would be interesting to see how they perform on the original data set without data cleaning. But this is subject of another research...

Logistic regression does require conversion of alphanumeric values to numeric, struggles dealing with the outliers and performs best when fitted with the data that have normal distribution.

Verdict (To be finalized) Despite sensitivity to data quality Logistic Regression outperforms other models in all other major categories.

Model Deployment

Without a doubt it would be a stretch to compare our model to the production numerical weather prediction models. But we do believe it might have a real live application as an educational tool. The model can demonstrate how various weather elements affect the probability of the rain.

It is simple to understand and deploy. The model does not require frequent updates because the weather patterns tend to be stable for a given geographical area (though this statement might be compromised in the context of the global warming). The model would benefit greatly if more complete data was available. Recall that we had to impute a lot of missing values.

Conclusion

Through exploring weather observations collected by 49 stations in Australia from 2007 to 2017 we selected and tuned a model to predict a rainy day tomorrow employing current day observations and historical data.

We commenced our research analyzing and understanding available data and geography of the weather stations. Then we identified the missing data, its distribution and feasibility of imputing it. We applied sophisticated data imputation algorithm to attack the problem. We continued our research selecting the most impactful data attributes to use as an input for our future model. Again we apply the feature identification algorithm to do the job.

When the data preparation phase was finished we picked and analysed four different classification models: Decision Tree, Naive Bayes, Random Forest and Logistic Regression. We conducted comparative analysis of the models, reviewed their strength and weaknesses. We fitted each model using K-fold cross-validation technique. Subsequently we evaluated performance of each model applying them to the test data set and comparing AUC - ROC and balanced accuracy metrics.

Finally we moved to identifying a winning model. In order to so we reviewed each model from different angles namely: * performance * interpretability * data quality sensitivity and data preparation effort The winning model scored the highest in the majority of the categories. It was Logistic Regression, which we employed to build a Shiny App Web application.

We consider the project to be a success. Being easily understood, with a balanced accuracy of 80% we conclude that our model could be applied for a short-term rain forecast. The last but not least we are confident that the model can predict the rain better than aboriginal "rainmakers".

Note from the Authors

This file was generated using [The R Journal style article template](#), additional information on how to prepare articles for submission is here - [Instructions for Authors](#). The article itself is an executable R Markdown file that could be [downloaded from Github](#) with all the necessary artifacts.

Sumaira Afzal
York University School of Continuing Studies

Viraja Ketkar
York University School of Continuing Studies

Murlidhar Loka
York University School of Continuing Studies

Vadim Spirkov
York University School of Continuing Studies