# Logistic Regression with Python and Scikit-Learn

In this project, I implement Logistic Regression with Python and Scikit-Learn. I build a classifier to predict whether or not it will rain tomorrow in Australia by training a binary classification model using Logistic Regression. I have used the **Rain in Australia** dataset downloaded from the Kaggle website for this project.

## Table of Contents

The table of contents for this project is as follows:-

# 1. Introduction to Logistic Regression

When data scientists may come across a new classification problem, the first algorithm that may come across their mind is **Logistic Regression**. It is a supervised learning classification algorithm which is used to predict observations to a discrete set of classes. Practically, it is used to classify observations into different categories. Hence, its output is discrete in nature. **Logistic Regression** is also called **Logit Regression**. It is one of the most simple, straightforward and versatile classification algorithms which is used to solve classification problems.

# 2. Logistic Regression intuition

In statistics, the **Logistic Regression model** is a widely used statistical model which is primarily used for classification purposes. It means that given a set of observations, Logistic Regression algorithm helps us to classify these observations into two or more discrete classes. So, the target variable is discrete in nature.

Logistic Regression algorithm works by implementing a linear equation with independent or explanatory variables to predict a response value. This predicted response value, denoted by z is then converted into a probability value that lie between 0 and 1. We use the **sigmoid function** in order to map predicted values to probability values. This sigmoid function then maps any real value into a probability value between 0 and 1.

The sigmoid function returns a probability value between 0 and 1. This probability value is then mapped to a discrete class which is either "0" or "1". In order to map this probability value to a discrete class (pass/fail, yes/no, true/false), we select a threshold value. This threshold value is called **Decision boundary**. Above this threshold value, we will map the probability values into class 1 and below which we will map values into class 0.

Mathematically, it can be expressed as follows:-

```
                p ≥ 0.5 => class = 1

                p < 0.5 => class = 0
```

Generally, the decision boundary is set to 0.5. So, if the probability value is 0.8 (> 0.5), we will map this observation to class 1. Similarly, if the probability value is 0.2 (< 0.5), we will map this observation to class 0.

We can use our knowledge of `sigmoid function` and `decision boundary` to write a prediction function. A prediction function in logistic regression returns the probability of the observation being positive, `Yes` or `True`. We call this as `class 1` and it is denoted by `P(class = 1)`. If the probability inches closer to one, then we will be more confident about our model that the observation is in class 1.

Logistic regression intuition is discussed in depth in the readme document.

# 3. The problem statement

In this project, I try to answer the question that whether or not it will rain tomorrow in Australia. I implement Logistic Regression with Python and Scikit-Learn.

To answer the question, I build a classifier to predict whether or not it will rain tomorrow in Australia by training a binary classification model using Logistic Regression. I have used the **Rain in Australia** dataset downloaded from the Kaggle website for this project.

# 4. Dataset description

I have used the **Rain in Australia** data set downloaded from the Kaggle website.

I have downloaded this data set from the Kaggle website. The data set can be found at the following url:-

https://www.kaggle.com/jsphyg/weather-dataset-rattle-package

This dataset contains daily weather observations from numerous Australian weather stations.

# 5. Import libraries

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

import warnings

warnings.filterwarnings('ignore')
```

# 6. Import dataset

```python
data = 'C:/datasets/weatherAUS.csv'

df = pd.read_csv(data)
```

# 7. Exploratory data analysis

Now, I will explore the data to gain insights about the data.

```python
# view dimensions of dataset

df.shape

(142193, 24)
```

We can see that there are 142193 instances and 24 variables in the data set.

```python
# preview the dataset

df.head()
```

```
         Date Location  MinTemp  MaxTemp  Rainfall  Evaporation
Sunshine  \
0  2008-12-01   Albury     13.4     22.9       0.6          NaN
NaN
1  2008-12-02   Albury      7.4     25.1       0.0          NaN
NaN
2  2008-12-03   Albury     12.9     25.7       0.0          NaN
NaN
3  2008-12-04   Albury      9.2     28.0       0.0          NaN
NaN
4  2008-12-05   Albury     17.5     32.3       1.0          NaN
NaN

  WindGustDir  WindGustSpeed WindDir9am    ...      Humidity3pm  \
0          W            44.0         W     ...             22.0
1         WNW           44.0        NNW    ...             25.0
2         WSW           46.0         W     ...             30.0
3          NE           24.0        SE     ...             16.0
4          W            41.0        ENE    ...             33.0

    Pressure9am  Pressure3pm  Cloud9am  Cloud3pm  Temp9am  Temp3pm
RainToday  \
0        1007.7       1007.1       8.0       NaN     16.9     21.8
No
1        1010.6       1007.8       NaN       NaN     17.2     24.3
No
2        1007.6       1008.7       NaN       2.0     21.0     23.2
No
3        1017.6       1012.8       NaN       NaN     18.1     26.5
No
4        1010.8       1006.0       7.0       8.0     17.8     29.7
No

    RISK_MM  RainTomorrow
0      0.0            No
1      0.0            No
2      0.0            No
3      1.0            No
4      0.2            No

[5 rows x 24 columns]

col_names = df.columns

col_names

Index(['Date', 'Location', 'MinTemp', 'MaxTemp', 'Rainfall',
'Evaporation',
       'Sunshine', 'WindGustDir', 'WindGustSpeed', 'WindDir9am',
'WindDir3pm',
```

```
         'WindSpeed9am', 'WindSpeed3pm', 'Humidity9am', 'Humidity3pm',
         'Pressure9am', 'Pressure3pm', 'Cloud9am', 'Cloud3pm',
'Temp9am',
         'Temp3pm', 'RainToday', 'RISK_MM', 'RainTomorrow'],
       dtype='object')
```

## Drop RISK_MM variable

It is given in the dataset description, that we should drop the RISK_MM feature variable from the dataset description. So, we should drop it as follows-

```
df.drop(['RISK_MM'], axis=1, inplace=True)

# view summary of dataset

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 142193 entries, 0 to 142192
Data columns (total 23 columns):
Date            142193 non-null object
Location        142193 non-null object
MinTemp         141556 non-null float64
MaxTemp         141871 non-null float64
Rainfall        140787 non-null float64
Evaporation      81350 non-null float64
Sunshine         74377 non-null float64
WindGustDir     132863 non-null object
WindGustSpeed   132923 non-null float64
WindDir9am      132180 non-null object
WindDir3pm      138415 non-null object
WindSpeed9am    140845 non-null float64
WindSpeed3pm    139563 non-null float64
Humidity9am     140419 non-null float64
Humidity3pm     138583 non-null float64
Pressure9am     128179 non-null float64
Pressure3pm     128212 non-null float64
Cloud9am         88536 non-null float64
Cloud3pm         85099 non-null float64
Temp9am         141289 non-null float64
Temp3pm         139467 non-null float64
RainToday       140787 non-null object
RainTomorrow    142193 non-null object
dtypes: float64(16), object(7)
memory usage: 25.0+ MB
```

## Types of variables

In this section, I segregate the dataset into categorical and numerical variables. There are a mixture of categorical and numerical variables in the dataset. Categorical variables have data type object. Numerical variables have data type float64.

First of all, I will find categorical variables.

```
# find categorical variables

categorical = [var for var in df.columns if df[var].dtype=='O']

print('There are {} categorical variables\n'.format(len(categorical)))

print('The categorical variables are :', categorical)

There are 7 categorical variables

The categorical variables are : ['Date', 'Location', 'WindGustDir',
'WindDir9am', 'WindDir3pm', 'RainToday', 'RainTomorrow']

# view the categorical variables

df[categorical].head()

        Date Location WindGustDir WindDir9am WindDir3pm RainToday  \
0  2008-12-01   Albury           W          W        WNW        No
1  2008-12-02   Albury         WNW        NNW        WSW        No
2  2008-12-03   Albury         WSW          W        WSW        No
3  2008-12-04   Albury          NE         SE          E        No
4  2008-12-05   Albury           W        ENE         NW        No

   RainTomorrow
0           No
1           No
2           No
3           No
4           No
```

## Summary of categorical variables

- There is a date variable. It is denoted by `Date` column.

- There are 6 categorical variables. These are given by `Location`, `WindGustDir`, `WindDir9am`, `WindDir3pm`, `RainToday` and `RainTomorrow`.

- There are two binary categorical variables - `RainToday` and `RainTomorrow`.

- `RainTomorrow` is the target variable.

# Explore problems within categorical variables

First, I will explore the categorical variables.

## Missing values in categorical variables

```python
# check missing values in categorical variables

df[categorical].isnull().sum()

Date                0
Location            0
WindGustDir      9330
WindDir9am      10013
WindDir3pm       3778
RainToday        1406
RainTomorrow        0
dtype: int64

# print categorical variables containing missing values

cat1 = [var for var in categorical if df[var].isnull().sum()!=0]

print(df[cat1].isnull().sum())

WindGustDir      9330
WindDir9am      10013
WindDir3pm       3778
RainToday        1406
dtype: int64
```

We can see that there are only 4 categorical variables in the dataset which contains missing values. These are `WindGustDir`, `WindDir9am`, `WindDir3pm` and `RainToday`.

## Frequency counts of categorical variables

Now, I will check the frequency counts of categorical variables.

```python
# view frequency of categorical variables

for var in categorical:

    print(df[var].value_counts())

2014-10-12    49
2017-01-15    49
2013-10-02    49
2014-07-15    49
2014-02-19    49
2016-08-21    49
2014-07-03    49
```

```
2016-10-21    49
2013-03-11    49
2017-02-08    49
2014-11-17    49
2013-04-25    49
2014-11-19    49
2014-08-30    49
2014-01-07    49
2013-04-10    49
2017-03-16    49
2013-09-04    49
2016-08-16    49
2016-10-19    49
2014-08-20    49
2017-05-12    49
2014-01-16    49
2016-07-22    49
2017-01-22    49
2013-09-25    49
2013-06-02    49
2016-07-06    49
2014-04-21    49
2013-10-16    49
              ..
2007-11-23     1
2008-01-15     1
2007-12-22     1
2007-11-08     1
2007-11-29     1
2008-01-29     1
2008-01-06     1
2007-11-02     1
2007-12-25     1
2008-01-28     1
2007-12-08     1
2007-11-09     1
2008-01-05     1
2007-11-26     1
2007-11-10     1
2007-11-20     1
2008-01-14     1
2007-12-03     1
2008-01-12     1
2007-11-03     1
2007-12-02     1
2008-01-31     1
2007-12-01     1
2007-11-06     1
2007-11-27     1
```

```
2007-12-19      1
2007-11-19      1
2007-12-30      1
2007-12-23      1
2008-01-09      1
Name: Date, Length: 3436, dtype: int64
Canberra             3418
Sydney               3337
Perth                3193
Darwin               3192
Hobart               3188
Brisbane             3161
Adelaide             3090
Bendigo              3034
Townsville           3033
AliceSprings         3031
MountGambier         3030
Launceston           3028
Ballarat             3028
Albany               3016
Albury               3011
PerthAirport         3009
MelbourneAirport     3009
Mildura              3007
SydneyAirport        3005
Nuriootpa            3002
Sale                 3000
Watsonia             2999
Tuggeranong          2998
Portland             2996
Woomera              2990
Cairns               2988
Cobar                2988
Wollongong           2983
GoldCoast            2980
WaggaWagga           2976
NorfolkIsland        2964
Penrith              2964
Newcastle            2955
SalmonGums           2955
CoffsHarbour         2953
Witchcliffe          2952
Richmond             2951
Dartmoor             2943
NorahHead            2929
BadgerysCreek        2928
MountGinini          2907
Moree                2854
Walpole              2819
```

```
PearceRAAF            2762
Williamtown           2553
Melbourne             2435
Nhil                  1569
Katherine             1559
Uluru                 1521
Name: Location, dtype: int64
W         9780
SE        9309
E         9071
N         9033
SSE       8993
S         8949
WSW       8901
SW        8797
SSW       8610
WNW       8066
NW        8003
ENE       7992
ESE       7305
NE        7060
NNW       6561
NNE       6433
Name: WindGustDir, dtype: int64
N        11393
SE        9162
E         9024
SSE       8966
NW        8552
S         8493
W         8260
SW        8237
NNE       7948
NNW       7840
ENE       7735
ESE       7558
NE        7527
SSW       7448
WNW       7194
WSW       6843
Name: WindDir9am, dtype: int64
SE       10663
W         9911
S         9598
WSW       9329
SW        9182
SSE       9142
N         8667
WNW       8656
```

```
NW         8468
ESE        8382
E          8342
NE         8164
SSW        8010
NNW        7733
ENE        7724
NNE        6444
Name: WindDir3pm, dtype: int64
No      109332
Yes      31455
Name: RainToday, dtype: int64
No      110316
Yes      31877
Name: RainTomorrow, dtype: int64
```

```python
# view frequency distribution of categorical variables

for var in categorical:

    print(df[var].value_counts()/np.float(len(df)))
```

```
2014-10-12    0.000345
2017-01-15    0.000345
2013-10-02    0.000345
2014-07-15    0.000345
2014-02-19    0.000345
2016-08-21    0.000345
2014-07-03    0.000345
2016-10-21    0.000345
2013-03-11    0.000345
2017-02-08    0.000345
2014-11-17    0.000345
2013-04-25    0.000345
2014-11-19    0.000345
2014-08-30    0.000345
2014-01-07    0.000345
2013-04-10    0.000345
2017-03-16    0.000345
2013-09-04    0.000345
2016-08-16    0.000345
2016-10-19    0.000345
2014-08-20    0.000345
2017-05-12    0.000345
2014-01-16    0.000345
2016-07-22    0.000345
2017-01-22    0.000345
2013-09-25    0.000345
2013-06-02    0.000345
2016-07-06    0.000345
```

```
2014-04-21    0.000345
2013-10-16    0.000345
                   ...
2007-11-23    0.000007
2008-01-15    0.000007
2007-12-22    0.000007
2007-11-08    0.000007
2007-11-29    0.000007
2008-01-29    0.000007
2008-01-06    0.000007
2007-11-02    0.000007
2007-12-25    0.000007
2008-01-28    0.000007
2007-12-08    0.000007
2007-11-09    0.000007
2008-01-05    0.000007
2007-11-26    0.000007
2007-11-10    0.000007
2007-11-20    0.000007
2008-01-14    0.000007
2007-12-03    0.000007
2008-01-12    0.000007
2007-11-03    0.000007
2007-12-02    0.000007
2008-01-31    0.000007
2007-12-01    0.000007
2007-11-06    0.000007
2007-11-27    0.000007
2007-12-19    0.000007
2007-11-19    0.000007
2007-12-30    0.000007
2007-12-23    0.000007
2008-01-09    0.000007
Name: Date, Length: 3436, dtype: float64
Canberra          0.024038
Sydney            0.023468
Perth             0.022455
Darwin            0.022448
Hobart            0.022420
Brisbane          0.022230
Adelaide          0.021731
Bendigo           0.021337
Townsville        0.021330
AliceSprings      0.021316
MountGambier      0.021309
Launceston        0.021295
Ballarat          0.021295
Albany            0.021211
Albury            0.021175
```

```
PerthAirport          0.021161
MelbourneAirport      0.021161
Mildura               0.021147
SydneyAirport         0.021133
Nuriootpa             0.021112
Sale                  0.021098
Watsonia              0.021091
Tuggeranong           0.021084
Portland              0.021070
Woomera               0.021028
Cairns                0.021014
Cobar                 0.021014
Wollongong            0.020979
GoldCoast             0.020957
WaggaWagga            0.020929
NorfolkIsland         0.020845
Penrith               0.020845
Newcastle             0.020782
SalmonGums            0.020782
CoffsHarbour          0.020768
Witchcliffe           0.020761
Richmond              0.020753
Dartmoor              0.020697
NorahHead             0.020599
BadgerysCreek         0.020592
MountGinini           0.020444
Moree                 0.020071
Walpole               0.019825
PearceRAAF            0.019424
Williamtown           0.017954
Melbourne             0.017125
Nhil                  0.011034
Katherine             0.010964
Uluru                 0.010697
Name: Location, dtype: float64
W       0.068780
SE      0.065467
E       0.063794
N       0.063526
SSE     0.063245
S       0.062936
WSW     0.062598
SW      0.061867
SSW     0.060552
WNW     0.056726
NW      0.056283
ENE     0.056205
ESE     0.051374
NE      0.049651
```

```
NNW     0.046142
NNE     0.045241
Name: WindGustDir, dtype: float64
N       0.080123
SE      0.064434
E       0.063463
SSE     0.063055
NW      0.060144
S       0.059729
W       0.058090
SW      0.057928
NNE     0.055896
NNW     0.055136
ENE     0.054398
ESE     0.053153
NE      0.052935
SSW     0.052380
WNW     0.050593
WSW     0.048125
Name: WindDir9am, dtype: float64
SE      0.074990
W       0.069701
S       0.067500
WSW     0.065608
SW      0.064574
SSE     0.064293
N       0.060952
WNW     0.060875
NW      0.059553
ESE     0.058948
E       0.058667
NE      0.057415
SSW     0.056332
NNW     0.054384
ENE     0.054321
NNE     0.045319
Name: WindDir3pm, dtype: float64
No      0.768899
Yes     0.221213
Name: RainToday, dtype: float64
No      0.775819
Yes     0.224181
Name: RainTomorrow, dtype: float64
```

## Number of labels: cardinality

The number of labels within a categorical variable is known as **cardinality**. A high number of labels within a variable is known as **high cardinality**. High cardinality may pose some serious problems in the machine learning model. So, I will check for high cardinality.

```python
# check for cardinality in categorical variables

for var in categorical:

    print(var, ' contains ', len(df[var].unique()), ' labels')
```

```
Date  contains  3436  labels
Location  contains  49  labels
WindGustDir  contains  17  labels
WindDir9am  contains  17  labels
WindDir3pm  contains  17  labels
RainToday  contains  3  labels
RainTomorrow  contains  2  labels
```

We can see that there is a `Date` variable which needs to be preprocessed. I will do preprocessing in the following section.

All the other variables contain relatively smaller number of variables.

## Feature Engineering of Date Variable

```python
df['Date'].dtypes
```

```
dtype('O')
```

We can see that the data type of `Date` variable is object. I will parse the date currently coded as object into datetime format.

```python
# parse the dates, currently coded as strings, into datetime format

df['Date'] = pd.to_datetime(df['Date'])

# extract year from date

df['Year'] = df['Date'].dt.year

df['Year'].head()
```

```
0    2008
1    2008
2    2008
3    2008
4    2008
Name: Year, dtype: int64
```

```python
# extract month from date

df['Month'] = df['Date'].dt.month

df['Month'].head()
```

```
0    12
1    12
2    12
3    12
4    12
Name: Month, dtype: int64
```

# extract day from date

```
df['Day'] = df['Date'].dt.day

df['Day'].head()

0    1
1    2
2    3
3    4
4    5
Name: Day, dtype: int64
```

# again view the summary of dataset

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 142193 entries, 0 to 142192
Data columns (total 26 columns):
Date             142193 non-null datetime64[ns]
Location         142193 non-null object
MinTemp          141556 non-null float64
MaxTemp          141871 non-null float64
Rainfall         140787 non-null float64
Evaporation      81350 non-null float64
Sunshine         74377 non-null float64
WindGustDir      132863 non-null object
WindGustSpeed    132923 non-null float64
WindDir9am       132180 non-null object
WindDir3pm       138415 non-null object
WindSpeed9am     140845 non-null float64
WindSpeed3pm     139563 non-null float64
Humidity9am      140419 non-null float64
Humidity3pm      138583 non-null float64
Pressure9am      128179 non-null float64
Pressure3pm      128212 non-null float64
Cloud9am         88536 non-null float64
Cloud3pm         85099 non-null float64
Temp9am          141289 non-null float64
Temp3pm          139467 non-null float64
RainToday        140787 non-null object
RainTomorrow     142193 non-null object
```

```
Year                142193 non-null int64
Month               142193 non-null int64
Day                 142193 non-null int64
dtypes: datetime64[ns](1), float64(16), int64(3), object(6)
memory usage: 28.2+ MB
```

We can see that there are three additional columns created from `Date` variable. Now, I will drop the original `Date` variable from the dataset.

```
# drop the original Date variable

df.drop('Date', axis=1, inplace = True)

# preview the dataset again

df.head()

  Location  MinTemp  MaxTemp  Rainfall  Evaporation  Sunshine
WindGustDir  \
0   Albury     13.4     22.9       0.6          NaN       NaN
W
1   Albury      7.4     25.1       0.0          NaN       NaN
WNW
2   Albury     12.9     25.7       0.0          NaN       NaN
WSW
3   Albury      9.2     28.0       0.0          NaN       NaN
NE
4   Albury     17.5     32.3       1.0          NaN       NaN
W

   WindGustSpeed WindDir9am WindDir3pm ...   Pressure3pm  Cloud9am
Cloud3pm  \
0            44.0          W        WNW ...        1007.1       8.0
NaN
1            44.0        NNW        WSW ...        1007.8       NaN
NaN
2            46.0          W        WSW ...        1008.7       NaN
2.0
3            24.0         SE          E ...        1012.8       NaN
NaN
4            41.0        ENE         NW ...        1006.0       7.0
8.0

   Temp9am  Temp3pm  RainToday  RainTomorrow  Year  Month  Day
0     16.9     21.8         No            No  2008     12    1
1     17.2     24.3         No            No  2008     12    2
2     21.0     23.2         No            No  2008     12    3
3     18.1     26.5         No            No  2008     12    4
4     17.8     29.7         No            No  2008     12    5
```

```
[5 rows x 25 columns]
```

Now, we can see that the `Date` variable has been removed from the dataset.

## Explore Categorical Variables

Now, I will explore the categorical variables one by one.

```python
# find categorical variables

categorical = [var for var in df.columns if df[var].dtype=='O']

print('There are {} categorical variables\n'.format(len(categorical)))

print('The categorical variables are :', categorical)
```
```
There are 6 categorical variables

The categorical variables are : ['Location', 'WindGustDir',
'WindDir9am', 'WindDir3pm', 'RainToday', 'RainTomorrow']
```

We can see that there are 6 categorical variables in the dataset. The `Date` variable has been removed. First, I will check missing values in categorical variables.

```python
# check for missing values in categorical variables

df[categorical].isnull().sum()
```
```
Location            0
WindGustDir      9330
WindDir9am      10013
WindDir3pm       3778
RainToday        1406
RainTomorrow        0
dtype: int64
```

We can see that `WindGustDir`, `WindDir9am`, `WindDir3pm`, `RainToday` variables contain missing values. I will explore these variables one by one.

## Explore `Location` variable

```python
# print number of labels in Location variable

print('Location contains', len(df.Location.unique()), 'labels')
```
```
Location contains 49 labels
```

```python
# check labels in location variable

df.Location.unique()

array(['Albury', 'BadgerysCreek', 'Cobar', 'CoffsHarbour', 'Moree',
       'Newcastle', 'NorahHead', 'NorfolkIsland', 'Penrith',
'Richmond',
       'Sydney', 'SydneyAirport', 'WaggaWagga', 'Williamtown',
       'Wollongong', 'Canberra', 'Tuggeranong', 'MountGinini',
'Ballarat',
       'Bendigo', 'Sale', 'MelbourneAirport', 'Melbourne', 'Mildura',
       'Nhil', 'Portland', 'Watsonia', 'Dartmoor', 'Brisbane',
'Cairns',
       'GoldCoast', 'Townsville', 'Adelaide', 'MountGambier',
'Nuriootpa',
       'Woomera', 'Albany', 'Witchcliffe', 'PearceRAAF',
'PerthAirport',
       'Perth', 'SalmonGums', 'Walpole', 'Hobart', 'Launceston',
       'AliceSprings', 'Darwin', 'Katherine', 'Uluru'], dtype=object)
```

```python
# check frequency distribution of values in Location variable

df.Location.value_counts()
```

```
Canberra          3418
Sydney            3337
Perth             3193
Darwin            3192
Hobart            3188
Brisbane          3161
Adelaide          3090
Bendigo           3034
Townsville        3033
AliceSprings      3031
MountGambier      3030
Launceston        3028
Ballarat          3028
Albany            3016
Albury            3011
PerthAirport      3009
MelbourneAirport  3009
Mildura           3007
SydneyAirport     3005
Nuriootpa         3002
Sale              3000
Watsonia          2999
Tuggeranong       2998
Portland          2996
Woomera           2990
Cairns            2988
```

```
Cobar                 2988
Wollongong            2983
GoldCoast             2980
WaggaWagga            2976
NorfolkIsland         2964
Penrith               2964
Newcastle             2955
SalmonGums            2955
CoffsHarbour          2953
Witchcliffe           2952
Richmond              2951
Dartmoor              2943
NorahHead             2929
BadgerysCreek         2928
MountGinini           2907
Moree                 2854
Walpole               2819
PearceRAAF            2762
Williamtown           2553
Melbourne             2435
Nhil                  1569
Katherine             1559
Uluru                 1521
Name: Location, dtype: int64
```

```python
# let's do One Hot Encoding of Location variable
# get k-1 dummy variables after One Hot Encoding
# preview the dataset with head() method

pd.get_dummies(df.Location, drop_first=True).head()
```

```
   Albany  Albury  AliceSprings  BadgerysCreek  Ballarat  Bendigo
Brisbane  \
0       0       1             0              0         0        0
0
1       0       1             0              0         0        0
0
2       0       1             0              0         0        0
0
3       0       1             0              0         0        0
0
4       0       1             0              0         0        0
0

   Cairns  Canberra  Cobar  ...  Townsville  Tuggeranong  Uluru  \
0       0         0      0  ...           0            0      0
1       0         0      0  ...           0            0      0
2       0         0      0  ...           0            0      0
3       0         0      0  ...           0            0      0
4       0         0      0  ...           0            0      0
```

```
     WaggaWagga  Walpole  Watsonia  Williamtown  Witchcliffe  Wollongong
\
0               0        0         0            0            0           0

1               0        0         0            0            0           0

2               0        0         0            0            0           0

3               0        0         0            0            0           0

4               0        0         0            0            0           0


     Woomera
0          0
1          0
2          0
3          0
4          0

[5 rows x 48 columns]
```

## Explore `WindGustDir` variable

```python
# print number of labels in WindGustDir variable

print('WindGustDir contains', len(df['WindGustDir'].unique()),
'labels')
```

```
WindGustDir contains 17 labels
```

```python
# check labels in WindGustDir variable

df['WindGustDir'].unique()
```

```
array(['W', 'WNW', 'WSW', 'NE', 'NNW', 'N', 'NNE', 'SW', 'ENE', 'SSE',
       'S', 'NW', 'SE', 'ESE', nan, 'E', 'SSW'], dtype=object)
```

```python
# check frequency distribution of values in WindGustDir variable

df.WindGustDir.value_counts()
```

```
W      9780
SE     9309
E      9071
N      9033
SSE    8993
S      8949
WSW    8901
SW     8797
```

```
SSW      8610
WNW      8066
NW       8003
ENE      7992
ESE      7305
NE       7060
NNW      6561
NNE      6433
Name: WindGustDir, dtype: int64
```

```python
# let's do One Hot Encoding of WindGustDir variable
# get k-1 dummy variables after One Hot Encoding
# also add an additional dummy variable to indicate there was missing data
# preview the dataset with head() method

pd.get_dummies(df.WindGustDir, drop_first=True, dummy_na=True).head()
```

| | ENE | ESE | N | NE | NNE | NNW | NW | S | SE | SSE | SSW | SW | W | WNW | WSW | NaN |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 3 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |

```python
# sum the number of 1s per boolean variable over the rows of the dataset
# it will tell us how many observations we have for each category

pd.get_dummies(df.WindGustDir, drop_first=True,
dummy_na=True).sum(axis=0)
```

```
ENE      7992
ESE      7305
N        9033
NE       7060
NNE      6433
NNW      6561
NW       8003
S        8949
SE       9309
SSE      8993
SSW      8610
SW       8797
```

```
W          9780
WNW        8066
WSW        8901
NaN        9330
dtype: int64
```

We can see that there are 9330 missing values in WindGustDir variable.

## Explore `WindDir9am` variable

```
# print number of labels in WindDir9am variable

print('WindDir9am contains', len(df['WindDir9am'].unique()), 'labels')

WindDir9am contains 17 labels

# check labels in WindDir9am variable

df['WindDir9am'].unique()

array(['W', 'NNW', 'SE', 'ENE', 'SW', 'SSE', 'S', 'NE', nan, 'SSW', 'N',
       'WSW', 'ESE', 'E', 'NW', 'WNW', 'NNE'], dtype=object)

# check frequency distribution of values in WindDir9am variable

df['WindDir9am'].value_counts()

N      11393
SE      9162
E       9024
SSE     8966
NW      8552
S       8493
W       8260
SW      8237
NNE     7948
NNW     7840
ENE     7735
ESE     7558
NE      7527
SSW     7448
WNW     7194
WSW     6843
Name: WindDir9am, dtype: int64

# let's do One Hot Encoding of WindDir9am variable
# get k-1 dummy variables after One Hot Encoding
# also add an additional dummy variable to indicate there was missing
data
# preview the dataset with head() method
```

```
pd.get_dummies(df.WindDir9am, drop_first=True, dummy_na=True).head()
```

| | ENE | ESE | N | NE | NNE | NNW | NW | S | SE | SSE | SSW | SW | W | WNW | WSW | NaN |
|---|-----|-----|---|----|-----|-----|----|---|----|-----|-----|----|---|-----|-----|-----|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

```
# sum the number of 1s per boolean variable over the rows of the
dataset
# it will tell us how many observations we have for each category

pd.get_dummies(df.WindDir9am, drop_first=True,
dummy_na=True).sum(axis=0)
```

```
ENE       7735
ESE       7558
N        11393
NE        7527
NNE       7948
NNW       7840
NW        8552
S         8493
SE        9162
SSE       8966
SSW       7448
SW        8237
W         8260
WNW       7194
WSW       6843
NaN      10013
dtype: int64
```

We can see that there are 10013 missing values in the `WindDir9am` variable.

## Explore `WindDir3pm` variable

```
# print number of labels in WindDir3pm variable

print('WindDir3pm contains', len(df['WindDir3pm'].unique()), 'labels')
```

```
WindDir3pm contains 17 labels
```

```python
# check labels in WindDir3pm variable

df['WindDir3pm'].unique()
```

```
array(['WNW', 'WSW', 'E', 'NW', 'W', 'SSE', 'ESE', 'ENE', 'NNW',
'SSW',
       'SW', 'SE', 'N', 'S', 'NNE', nan, 'NE'], dtype=object)
```

```python
# check frequency distribution of values in WindDir3pm variable

df['WindDir3pm'].value_counts()
```

```
SE      10663
W        9911
S        9598
WSW      9329
SW       9182
SSE      9142
N        8667
WNW      8656
NW       8468
ESE      8382
E        8342
NE       8164
SSW      8010
NNW      7733
ENE      7724
NNE      6444
Name: WindDir3pm, dtype: int64
```

```python
# let's do One Hot Encoding of WindDir3pm variable
# get k-1 dummy variables after One Hot Encoding
# also add an additional dummy variable to indicate there was missing data
# preview the dataset with head() method

pd.get_dummies(df.WindDir3pm, drop_first=True, dummy_na=True).head()
```

|   | ENE | ESE | N | NE | NNE | NNW | NW | S | SE | SSE | SSW | SW | W | WNW | WSW | NaN |
|---|-----|-----|---|----|-----|-----|----|---|----|-----|-----|----|---|-----|-----|-----|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

```
# sum the number of 1s per boolean variable over the rows of the
dataset
# it will tell us how many observations we have for each category

pd.get_dummies(df.WindDir3pm, drop_first=True,
dummy_na=True).sum(axis=0)

ENE       7724
ESE       8382
N         8667
NE        8164
NNE       6444
NNW       7733
NW        8468
S         9598
SE       10663
SSE       9142
SSW       8010
SW        9182
W         9911
WNW       8656
WSW       9329
NaN       3778
dtype: int64
```

There are 3778 missing values in the `WindDir3pm` variable.

## Explore `RainToday` variable

```
# print number of labels in RainToday variable

print('RainToday contains', len(df['RainToday'].unique()), 'labels')

RainToday contains 3 labels

# check labels in WindGustDir variable

df['RainToday'].unique()

array(['No', 'Yes', nan], dtype=object)

# check frequency distribution of values in WindGustDir variable

df.RainToday.value_counts()

No       109332
Yes       31455
Name: RainToday, dtype: int64

# let's do One Hot Encoding of RainToday variable
# get k-1 dummy variables after One Hot Encoding
```

```
# also add an additional dummy variable to indicate there was missing
data
# preview the dataset with head() method

pd.get_dummies(df.RainToday, drop_first=True, dummy_na=True).head()

    Yes  NaN
0    0    0
1    0    0
2    0    0
3    0    0
4    0    0

# sum the number of 1s per boolean variable over the rows of the
dataset
# it will tell us how many observations we have for each category

pd.get_dummies(df.RainToday, drop_first=True,
dummy_na=True).sum(axis=0)

Yes    31455
NaN     1406
dtype: int64
```

There are 1406 missing values in the `RainToday` variable.

## Explore Numerical Variables

```
# find numerical variables

numerical = [var for var in df.columns if df[var].dtype!='O']

print('There are {} numerical variables\n'.format(len(numerical)))

print('The numerical variables are :', numerical)

There are 19 numerical variables

The numerical variables are : ['MinTemp', 'MaxTemp', 'Rainfall',
'Evaporation', 'Sunshine', 'WindGustSpeed', 'WindSpeed9am',
'WindSpeed3pm', 'Humidity9am', 'Humidity3pm', 'Pressure9am',
'Pressure3pm', 'Cloud9am', 'Cloud3pm', 'Temp9am', 'Temp3pm', 'Year',
'Month', 'Day']

# view the numerical variables

df[numerical].head()

   MinTemp  MaxTemp  Rainfall  Evaporation  Sunshine  WindGustSpeed  \
0     13.4     22.9       0.6          NaN       NaN           44.0
1      7.4     25.1       0.0          NaN       NaN           44.0
```

|   | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine | WindGustSpeed |
|---|---|---|---|---|---|---|
| 2 | 12.9 | 25.7 | 0.0 | NaN | NaN | 46.0 |
| 3 | 9.2 | 28.0 | 0.0 | NaN | NaN | 24.0 |
| 4 | 17.5 | 32.3 | 1.0 | NaN | NaN | 41.0 |

|   | WindSpeed9am | WindSpeed3pm | Humidity9am | Humidity3pm | Pressure9am \ |
|---|---|---|---|---|---|
| 0 | 20.0 | 24.0 | 71.0 | 22.0 | 1007.7 |
| 1 | 4.0 | 22.0 | 44.0 | 25.0 | 1010.6 |
| 2 | 19.0 | 26.0 | 38.0 | 30.0 | 1007.6 |
| 3 | 11.0 | 9.0 | 45.0 | 16.0 | 1017.6 |
| 4 | 7.0 | 20.0 | 82.0 | 33.0 | 1010.8 |

|   | Pressure3pm | Cloud9am | Cloud3pm | Temp9am | Temp3pm | Year | Month | Day |
|---|---|---|---|---|---|---|---|---|
| 0 | 1007.1 | 8.0 | NaN | 16.9 | 21.8 | 2008 | 12 | 1 |
| 1 | 1007.8 | NaN | NaN | 17.2 | 24.3 | 2008 | 12 | 2 |
| 2 | 1008.7 | NaN | 2.0 | 21.0 | 23.2 | 2008 | 12 | 3 |
| 3 | 1012.8 | NaN | NaN | 18.1 | 26.5 | 2008 | 12 | 4 |
| 4 | 1006.0 | 7.0 | 8.0 | 17.8 | 29.7 | 2008 | 12 | 5 |

## Summary of numerical variables

- There are 16 numerical variables.

- These are given by `MinTemp`, `MaxTemp`, `Rainfall`, `Evaporation`, `Sunshine`, `WindGustSpeed`, `WindSpeed9am`, `WindSpeed3pm`, `Humidity9am`, `Humidity3pm`, `Pressure9am`, `Pressure3pm`, `Cloud9am`, `Cloud3pm`, `Temp9am` and `Temp3pm`.

- All of the numerical variables are of continuous type.

# Explore problems within numerical variables

Now, I will explore the numerical variables.

## Missing values in numerical variables

```
# check missing values in numerical variables

df[numerical].isnull().sum()
```

```
MinTemp            637
MaxTemp            322
```

```
Rainfall             1406
Evaporation         60843
Sunshine            67816
WindGustSpeed        9270
WindSpeed9am         1348
WindSpeed3pm         2630
Humidity9am          1774
Humidity3pm          3610
Pressure9am         14014
Pressure3pm         13981
Cloud9am            53657
Cloud3pm            57094
Temp9am               904
Temp3pm              2726
Year                    0
Month                   0
Day                     0
dtype: int64
```

We can see that all the 16 numerical variables contain missing values.

## Outliers in numerical variables

```python
# view summary statistics in numerical variables

print(round(df[numerical].describe()),2)
```

```
        MinTemp   MaxTemp  Rainfall  Evaporation  Sunshine
WindGustSpeed  \
count  141556.0  141871.0  140787.0      81350.0   74377.0
132923.0
mean       12.0      23.0       2.0          5.0       8.0
40.0
std         6.0       7.0       8.0          4.0       4.0
14.0
min        -8.0      -5.0       0.0          0.0       0.0
6.0
25%         8.0      18.0       0.0          3.0       5.0
31.0
50%        12.0      23.0       0.0          5.0       8.0
39.0
75%        17.0      28.0       1.0          7.0      11.0
48.0
max        34.0      48.0     371.0        145.0      14.0
135.0

        WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm
Pressure9am  \
count       140845.0      139563.0     140419.0      138583.0
128179.0
```

|       |       |       |       |       |        |
|-------|-------|-------|-------|-------|--------|
| mean  | 14.0  | 19.0  | 69.0  | 51.0  | 1018.0 |
| std   | 9.0   | 9.0   | 19.0  | 21.0  | 7.0    |
| min   | 0.0   | 0.0   | 0.0   | 0.0   | 980.0  |
| 25%   | 7.0   | 13.0  | 57.0  | 37.0  | 1013.0 |
| 50%   | 13.0  | 19.0  | 70.0  | 52.0  | 1018.0 |
| 75%   | 19.0  | 24.0  | 83.0  | 66.0  | 1022.0 |
| max   | 130.0 | 87.0  | 100.0 | 100.0 | 1041.0 |

|       | Pressure3pm | Cloud9am | Cloud3pm | Temp9am  | Temp3pm  | Year     |
|-------|-------------|----------|----------|----------|----------|----------|
| count | 128212.0    | 88536.0  | 85099.0  | 141289.0 | 139467.0 | 142193.0 |
| mean  | 1015.0      | 4.0      | 5.0      | 17.0     | 22.0     | 2013.0   |
| std   | 7.0         | 3.0      | 3.0      | 6.0      | 7.0      | 3.0      |
| min   | 977.0       | 0.0      | 0.0      | -7.0     | -5.0     | 2007.0   |
| 25%   | 1010.0      | 1.0      | 2.0      | 12.0     | 17.0     | 2011.0   |
| 50%   | 1015.0      | 5.0      | 5.0      | 17.0     | 21.0     | 2013.0   |
| 75%   | 1020.0      | 7.0      | 7.0      | 22.0     | 26.0     | 2015.0   |
| max   | 1040.0      | 9.0      | 9.0      | 40.0     | 47.0     | 2017.0   |

|       | Month    | Day      |   |
|-------|----------|----------|---|
| count | 142193.0 | 142193.0 |   |
| mean  | 6.0      | 16.0     |   |
| std   | 3.0      | 9.0      |   |
| min   | 1.0      | 1.0      |   |
| 25%   | 3.0      | 8.0      |   |
| 50%   | 6.0      | 16.0     |   |
| 75%   | 9.0      | 23.0     |   |
| max   | 12.0     | 31.0     | 2 |

On closer inspection, we can see that the Rainfall, Evaporation, WindSpeed9am and WindSpeed3pm columns may contain outliers.

I will draw boxplots to visualise outliers in the above variables.

```python
# draw boxplots to visualize outliers

plt.figure(figsize=(15,10))


plt.subplot(2, 2, 1)
fig = df.boxplot(column='Rainfall')
fig.set_title('')
fig.set_ylabel('Rainfall')


plt.subplot(2, 2, 2)
fig = df.boxplot(column='Evaporation')
fig.set_title('')
fig.set_ylabel('Evaporation')


plt.subplot(2, 2, 3)
fig = df.boxplot(column='WindSpeed9am')
fig.set_title('')
fig.set_ylabel('WindSpeed9am')


plt.subplot(2, 2, 4)
fig = df.boxplot(column='WindSpeed3pm')
fig.set_title('')
fig.set_ylabel('WindSpeed3pm')

Text(0,0.5,'WindSpeed3pm')
```

The above boxplots confirm that there are lot of outliers in these variables.

## Check the distribution of variables

Now, I will plot the histograms to check distributions to find out if they are normal or skewed. If the variable follows normal distribution, then I will do `Extreme Value Analysis` otherwise if they are skewed, I will find IQR (Interquantile range).

```python
# plot histogram to check distribution

plt.figure(figsize=(15,10))


plt.subplot(2, 2, 1)
fig = df.Rainfall.hist(bins=10)
fig.set_xlabel('Rainfall')
fig.set_ylabel('RainTomorrow')


plt.subplot(2, 2, 2)
fig = df.Evaporation.hist(bins=10)
fig.set_xlabel('Evaporation')
fig.set_ylabel('RainTomorrow')


plt.subplot(2, 2, 3)
```

```
fig = df.WindSpeed9am.hist(bins=10)
fig.set_xlabel('WindSpeed9am')
fig.set_ylabel('RainTomorrow')


plt.subplot(2, 2, 4)
fig = df.WindSpeed3pm.hist(bins=10)
fig.set_xlabel('WindSpeed3pm')
fig.set_ylabel('RainTomorrow')

Text(0,0.5,'RainTomorrow')
```



We can see that all the four variables are skewed. So, I will use interquantile range to find outliers.

```
# find outliers for Rainfall variable

IQR = df.Rainfall.quantile(0.75) - df.Rainfall.quantile(0.25)
Lower_fence = df.Rainfall.quantile(0.25) - (IQR * 3)
Upper_fence = df.Rainfall.quantile(0.75) + (IQR * 3)
print('Rainfall outliers are values < {lowerboundary} or >
{upperboundary}'.format(lowerboundary=Lower_fence,
upperboundary=Upper_fence))

Rainfall outliers are values < -2.4000000000000004 or > 3.2
```

For `Rainfall`, the minimum and maximum values are 0.0 and 371.0. So, the outliers are values > 3.2.

```python
# find outliers for Evaporation variable

IQR = df.Evaporation.quantile(0.75) - df.Evaporation.quantile(0.25)
Lower_fence = df.Evaporation.quantile(0.25) - (IQR * 3)
Upper_fence = df.Evaporation.quantile(0.75) + (IQR * 3)
print('Evaporation outliers are values < {lowerboundary} or >
{upperboundary}'.format(lowerboundary=Lower_fence,
upperboundary=Upper_fence))

Evaporation outliers are values < -11.800000000000002 or >
21.800000000000004
```

For `Evaporation`, the minimum and maximum values are 0.0 and 145.0. So, the outliers are values > 21.8.

```python
# find outliers for WindSpeed9am variable

IQR = df.WindSpeed9am.quantile(0.75) - df.WindSpeed9am.quantile(0.25)
Lower_fence = df.WindSpeed9am.quantile(0.25) - (IQR * 3)
Upper_fence = df.WindSpeed9am.quantile(0.75) + (IQR * 3)
print('WindSpeed9am outliers are values < {lowerboundary} or >
{upperboundary}'.format(lowerboundary=Lower_fence,
upperboundary=Upper_fence))

WindSpeed9am outliers are values < -29.0 or > 55.0
```

For `WindSpeed9am`, the minimum and maximum values are 0.0 and 130.0. So, the outliers are values > 55.0.

```python
# find outliers for WindSpeed3pm variable

IQR = df.WindSpeed3pm.quantile(0.75) - df.WindSpeed3pm.quantile(0.25)
Lower_fence = df.WindSpeed3pm.quantile(0.25) - (IQR * 3)
Upper_fence = df.WindSpeed3pm.quantile(0.75) + (IQR * 3)
print('WindSpeed3pm outliers are values < {lowerboundary} or >
{upperboundary}'.format(lowerboundary=Lower_fence,
upperboundary=Upper_fence))

WindSpeed3pm outliers are values < -20.0 or > 57.0
```

For `WindSpeed3pm`, the minimum and maximum values are 0.0 and 87.0. So, the outliers are values > 57.0.

# 8. Declare feature vector and target variable

```
X = df.drop(['RainTomorrow'], axis=1)

y = df['RainTomorrow']
```

# 9. Split data into separate training and test set

```
# split X and y into training and testing sets

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
0.2, random_state = 0)

# check the shape of X_train and X_test

X_train.shape, X_test.shape

((113754, 24), (28439, 24))
```

# 10. Feature Engineering

**Feature Engineering** is the process of transforming raw data into useful features that help us to understand our model better and increase its predictive power. I will carry out feature engineering on different types of variables.

First, I will display the categorical and numerical variables again separately.

```
# check data types in X_train

X_train.dtypes

Location          object
MinTemp           float64
MaxTemp           float64
Rainfall          float64
Evaporation       float64
Sunshine          float64
WindGustDir       object
WindGustSpeed     float64
WindDir9am        object
WindDir3pm        object
WindSpeed9am      float64
WindSpeed3pm      float64
Humidity9am       float64
Humidity3pm       float64
Pressure9am       float64
Pressure3pm       float64
Cloud9am          float64
```

```
Cloud3pm         float64
Temp9am          float64
Temp3pm          float64
RainToday         object
Year               int64
Month              int64
Day                int64
dtype: object
```

```python
# display categorical variables

categorical = [col for col in X_train.columns if X_train[col].dtypes == 'O']

categorical
```

```
['Location', 'WindGustDir', 'WindDir9am', 'WindDir3pm', 'RainToday']
```

```python
# display numerical variables

numerical = [col for col in X_train.columns if X_train[col].dtypes != 'O']

numerical
```

```
['MinTemp',
 'MaxTemp',
 'Rainfall',
 'Evaporation',
 'Sunshine',
 'WindGustSpeed',
 'WindSpeed9am',
 'WindSpeed3pm',
 'Humidity9am',
 'Humidity3pm',
 'Pressure9am',
 'Pressure3pm',
 'Cloud9am',
 'Cloud3pm',
 'Temp9am',
 'Temp3pm',
 'Year',
 'Month',
 'Day']
```

## Engineering missing values in numerical variables

```python
# check missing values in numerical variables in X_train

X_train[numerical].isnull().sum()
```

```
MinTemp              495
MaxTemp              264
Rainfall            1139
Evaporation        48718
Sunshine           54314
WindGustSpeed       7367
WindSpeed9am        1086
WindSpeed3pm        2094
Humidity9am         1449
Humidity3pm         2890
Pressure9am        11212
Pressure3pm        11186
Cloud9am           43137
Cloud3pm           45768
Temp9am              740
Temp3pm             2171
Year                   0
Month                  0
Day                    0
dtype: int64
```

```python
# check missing values in numerical variables in X_test

X_test[numerical].isnull().sum()
```

```
MinTemp              142
MaxTemp               58
Rainfall             267
Evaporation        12125
Sunshine           13502
WindGustSpeed       1903
WindSpeed9am         262
WindSpeed3pm         536
Humidity9am          325
Humidity3pm          720
Pressure9am         2802
Pressure3pm         2795
Cloud9am           10520
Cloud3pm           11326
Temp9am              164
Temp3pm              555
Year                   0
Month                  0
Day                    0
dtype: int64
```

```python
# print percentage of missing values in the numerical variables in
training set

for col in numerical:
```

```
    if X_train[col].isnull().mean()>0:
        print(col, round(X_train[col].isnull().mean(),4))
```

```
MinTemp 0.0044
MaxTemp 0.0023
Rainfall 0.01
Evaporation 0.4283
Sunshine 0.4775
WindGustSpeed 0.0648
WindSpeed9am 0.0095
WindSpeed3pm 0.0184
Humidity9am 0.0127
Humidity3pm 0.0254
Pressure9am 0.0986
Pressure3pm 0.0983
Cloud9am 0.3792
Cloud3pm 0.4023
Temp9am 0.0065
Temp3pm 0.0191
```

## Assumption

I assume that the data are missing completely at random (MCAR). There are two methods which can be used to impute missing values. One is mean or median imputation and other one is random sample imputation. When there are outliers in the dataset, we should use median imputation. So, I will use median imputation because median imputation is robust to outliers.

I will impute missing values with the appropriate statistical measures of the data, in this case median. Imputation should be done over the training set, and then propagated to the test set. It means that the statistical measures to be used to fill missing values both in train and test set, should be extracted from the train set only. This is to avoid overfitting.

```python
# impute missing values in X_train and X_test with respective column
median in X_train

for df1 in [X_train, X_test]:
    for col in numerical:
        col_median=X_train[col].median()
        df1[col].fillna(col_median, inplace=True)


# check again missing values in numerical variables in X_train

X_train[numerical].isnull().sum()
```

```
MinTemp          0
MaxTemp          0
Rainfall         0
Evaporation      0
Sunshine         0
```

```
WindGustSpeed     0
WindSpeed9am      0
WindSpeed3pm      0
Humidity9am       0
Humidity3pm       0
Pressure9am       0
Pressure3pm       0
Cloud9am          0
Cloud3pm          0
Temp9am           0
Temp3pm           0
Year              0
Month             0
Day               0
dtype: int64
```

```python
# check missing values in numerical variables in X_test

X_test[numerical].isnull().sum()
```

```
MinTemp           0
MaxTemp           0
Rainfall          0
Evaporation       0
Sunshine          0
WindGustSpeed     0
WindSpeed9am      0
WindSpeed3pm      0
Humidity9am       0
Humidity3pm       0
Pressure9am       0
Pressure3pm       0
Cloud9am          0
Cloud3pm          0
Temp9am           0
Temp3pm           0
Year              0
Month             0
Day               0
dtype: int64
```

Now, we can see that there are no missing values in the numerical columns of training and test set.

## Engineering missing values in categorical variables

```python
# print percentage of missing values in the categorical variables in
# training set

X_train[categorical].isnull().mean()
```

```
Location        0.000000
WindGustDir     0.065114
WindDir9am      0.070134
WindDir3pm      0.026443
RainToday       0.010013
dtype: float64
```

```python
# print categorical variables with missing data

for col in categorical:
    if X_train[col].isnull().mean()>0:
        print(col, (X_train[col].isnull().mean()))
```

```
WindGustDir 0.06511419378659213
WindDir9am 0.07013379749283542
WindDir3pm 0.026443026179299188
RainToday 0.01001283471350458
```

```python
# impute missing categorical variables with most frequent value

for df2 in [X_train, X_test]:
    df2['WindGustDir'].fillna(X_train['WindGustDir'].mode()[0],
inplace=True)
    df2['WindDir9am'].fillna(X_train['WindDir9am'].mode()[0],
inplace=True)
    df2['WindDir3pm'].fillna(X_train['WindDir3pm'].mode()[0],
inplace=True)
    df2['RainToday'].fillna(X_train['RainToday'].mode()[0],
inplace=True)
```

```python
# check missing values in categorical variables in X_train

X_train[categorical].isnull().sum()
```

```
Location        0
WindGustDir     0
WindDir9am      0
WindDir3pm      0
RainToday       0
dtype: int64
```

```python
# check missing values in categorical variables in X_test

X_test[categorical].isnull().sum()
```

```
Location        0
WindGustDir     0
WindDir9am      0
WindDir3pm      0
RainToday       0
dtype: int64
```

As a final check, I will check for missing values in X_train and X_test.

```
# check missing values in X_train

X_train.isnull().sum()
```

```
Location          0
MinTemp           0
MaxTemp           0
Rainfall          0
Evaporation       0
Sunshine          0
WindGustDir       0
WindGustSpeed     0
WindDir9am        0
WindDir3pm        0
WindSpeed9am      0
WindSpeed3pm      0
Humidity9am       0
Humidity3pm       0
Pressure9am       0
Pressure3pm       0
Cloud9am          0
Cloud3pm          0
Temp9am           0
Temp3pm           0
RainToday         0
Year              0
Month             0
Day               0
dtype: int64
```

```
# check missing values in X_test

X_test.isnull().sum()
```

```
Location          0
MinTemp           0
MaxTemp           0
Rainfall          0
Evaporation       0
Sunshine          0
WindGustDir       0
WindGustSpeed     0
WindDir9am        0
WindDir3pm        0
WindSpeed9am      0
WindSpeed3pm      0
Humidity9am       0
Humidity3pm       0
```

```
Pressure9am        0
Pressure3pm        0
Cloud9am           0
Cloud3pm           0
Temp9am            0
Temp3pm            0
RainToday          0
Year               0
Month              0
Day                0
dtype: int64
```

We can see that there are no missing values in X_train and X_test.

## Engineering outliers in numerical variables

We have seen that the `Rainfall`, `Evaporation`, `WindSpeed9am` and `WindSpeed3pm` columns contain outliers. I will use top-coding approach to cap maximum values and remove outliers from the above variables.

```python
def max_value(df3, variable, top):
    return np.where(df3[variable]>top, top, df3[variable])

for df3 in [X_train, X_test]:
    df3['Rainfall'] = max_value(df3, 'Rainfall', 3.2)
    df3['Evaporation'] = max_value(df3, 'Evaporation', 21.8)
    df3['WindSpeed9am'] = max_value(df3, 'WindSpeed9am', 55)
    df3['WindSpeed3pm'] = max_value(df3, 'WindSpeed3pm', 57)

X_train.Rainfall.max(), X_test.Rainfall.max()

(3.2, 3.2)

X_train.Evaporation.max(), X_test.Evaporation.max()

(21.8, 21.8)

X_train.WindSpeed9am.max(), X_test.WindSpeed9am.max()

(55.0, 55.0)

X_train.WindSpeed3pm.max(), X_test.WindSpeed3pm.max()

(57.0, 57.0)

X_train[numerical].describe()
```

|       | MinTemp | MaxTemp | Rainfall | Evaporation \ |
|-------|---------|---------|----------|---------------|
| count | 113754.000000 | 113754.000000 | 113754.000000 | 113754.000000 |
| mean  | 12.193497 | 23.237216 | 0.675080 | 5.151606 |
| std   | 6.388279 | 7.094149 | 1.183837 | 2.823707 |

```
min          -8.200000        -4.800000         0.000000          0.000000
25%           7.600000        18.000000         0.000000          4.000000
50%          12.000000        22.600000         0.000000          4.800000
75%          16.800000        28.200000         0.600000          5.400000
max          33.900000        48.100000         3.200000         21.800000

              Sunshine    WindGustSpeed     WindSpeed9am     WindSpeed3pm   \
count    113754.000000    113754.000000    113754.000000    113754.000000
mean          8.041154        39.884074        13.978155        18.614756
std           2.769480        13.116959         8.806558         8.685862
min           0.000000         6.000000         0.000000         0.000000
25%           8.200000        31.000000         7.000000        13.000000
50%           8.500000        39.000000        13.000000        19.000000
75%           8.700000        46.000000        19.000000        24.000000
max          14.500000       135.000000        55.000000        57.000000

             Humidity9am      Humidity3pm       Pressure9am      Pressure3pm   \
count    113754.000000    113754.000000    113754.000000    113754.000000
mean         68.867486        51.509547      1017.640649      1015.241101
std          18.935587        20.530723         6.738680         6.675168
min           0.000000         0.000000       980.500000       977.100000
25%          57.000000        37.000000      1013.500000      1011.000000
50%          70.000000        52.000000      1017.600000      1015.200000
75%          83.000000        65.000000      1021.800000      1019.400000
max         100.000000       100.000000      1041.000000      1039.600000

                Cloud9am         Cloud3pm          Temp9am          Temp3pm   \
count    113754.000000    113754.000000    113754.000000    113754.000000
mean          4.651801         4.703588        16.995062        21.688643
std           2.292726         2.117847         6.463772         6.855649
min           0.000000         0.000000        -7.200000        -5.400000
25%           3.000000         4.000000        12.300000        16.700000
50%           5.000000         5.000000        16.700000        21.100000
75%           6.000000         6.000000        21.500000        26.300000
max           9.000000         8.000000        40.200000        46.700000

                    Year            Month              Day
count    113754.000000    113754.000000    113754.000000
mean       2012.759727         6.404021        15.710419
std           2.540419         3.427798         8.796821
min        2007.000000         1.000000         1.000000
25%        2011.000000         3.000000         8.000000
50%        2013.000000         6.000000        16.000000
75%        2015.000000         9.000000        23.000000
max        2017.000000        12.000000        31.000000
```

We can now see that the outliers in Rainfall, Evaporation, WindSpeed9am and WindSpeed3pm columns are capped.

## Encode categorical variables

```
categorical

['Location', 'WindGustDir', 'WindDir9am', 'WindDir3pm', 'RainToday']

X_train[categorical].head()
```

```
            Location WindGustDir WindDir9am WindDir3pm RainToday
110803   Witchcliffe           S        SSE          S        No
87289         Cairns         ENE        SSE         SE       Yes
134949   AliceSprings          E         NE          N        No
85553         Cairns         ESE        SSE          E        No
16110      Newcastle           W          N         SE        No
```

```python
# encode RainToday variable

import category_encoders as ce

encoder = ce.BinaryEncoder(cols=['RainToday'])

X_train = encoder.fit_transform(X_train)

X_test = encoder.transform(X_test)

X_train.head()
```

```
         RainToday_0  RainToday_1        Location  MinTemp  MaxTemp
Rainfall  \
110803             0            1     Witchcliffe     13.9     22.6
0.2
87289              1            0          Cairns     22.4     29.4
2.0
134949             0            1    AliceSprings      9.7     36.2
0.0
85553              0            1          Cairns     20.5     30.1
0.0
16110              0            1       Newcastle     16.8     29.2
0.0


         Evaporation  Sunshine WindGustDir  WindGustSpeed ...
Humidity3pm  \
110803           4.8       8.5           S           41.0 ...
55.0
87289            6.0       6.3         ENE           33.0 ...
59.0
134949          11.4      12.3           E           31.0 ...
2.0
85553            8.8      11.1         ESE           37.0 ...
53.0
16110            4.8       8.5           W           39.0 ...
53.0
```

|        | Pressure9am | Pressure3pm | Cloud9am | Cloud3pm | Temp9am | Temp3pm | Year |
|--------|-------------|-------------|----------|----------|---------|---------|------|
| 110803 | 1013.9      | 1013.4      | 5.0      | 5.0      | 18.8    | 20.4    | 2014 |
| 87289  | 1016.9      | 1013.1      | 7.0      | 5.0      | 26.4    | 27.5    | 2015 |
| 134949 | 1018.1      | 1013.6      | 1.0      | 1.0      | 28.5    | 35.0    | 2014 |
| 85553  | 1014.1      | 1010.8      | 2.0      | 3.0      | 27.3    | 29.4    | 2010 |
| 16110  | 1017.6      | 1015.2      | 5.0      | 8.0      | 22.2    | 27.0    | 2012 |

|        | Month | Day |
|--------|-------|-----|
| 110803 | 4     | 25  |
| 87289  | 11    | 2   |
| 134949 | 10    | 19  |
| 85553  | 10    | 30  |
| 16110  | 11    | 8   |

[5 rows x 25 columns]

We can see that two additional variables `RainToday_0` and `RainToday_1` are created from `RainToday` variable.

Now, I will create the `X_train` training set.

```
X_train = pd.concat([X_train[numerical], X_train[['RainToday_0',
'RainToday_1']],
                     pd.get_dummies(X_train.Location),
                     pd.get_dummies(X_train.WindGustDir),
                     pd.get_dummies(X_train.WindDir9am),
                     pd.get_dummies(X_train.WindDir3pm)], axis=1)

X_train.head()
```

|        | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine | WindGustSpeed |
|--------|---------|---------|----------|-------------|----------|---------------|
| 110803 | 13.9    | 22.6    | 0.2      | 4.8         | 8.5      | 41.0          |
| 87289  | 22.4    | 29.4    | 2.0      | 6.0         | 6.3      | 33.0          |
| 134949 | 9.7     | 36.2    | 0.0      | 11.4        | 12.3     | 31.0          |
| 85553  | 20.5    | 30.1    | 0.0      | 8.8         | 11.1     | 37.0          |
| 16110  | 16.8    | 29.2    | 0.0      | 4.8         | 8.5      | 39.0          |

```
        WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm ...    NNW
NW   S  \
110803          20.0          28.0         65.0         55.0 ...      0
0   1
87289            7.0          19.0         71.0         59.0 ...      0
0   0
134949          15.0          11.0          6.0          2.0 ...      0
0   0
85553           22.0          19.0         59.0         53.0 ...      0
0   0
16110            0.0           7.0         72.0         53.0 ...      0
0   0

        SE  SSE  SSW  SW  W  WNW  WSW
110803   0    0    0   0  0    0    0
87289    1    0    0   0  0    0    0
134949   0    0    0   0  0    0    0
85553    0    0    0   0  0    0    0
16110    1    0    0   0  0    0    0

[5 rows x 118 columns]
```

Similarly, I will create the `X_test` testing set.

```
X_test = pd.concat([X_test[numerical], X_test[['RainToday_0',
'RainToday_1']],
                    pd.get_dummies(X_test.Location),
                    pd.get_dummies(X_test.WindGustDir),
                    pd.get_dummies(X_test.WindDir9am),
                    pd.get_dummies(X_test.WindDir3pm)], axis=1)

X_test.head()

        MinTemp  MaxTemp  Rainfall  Evaporation  Sunshine
WindGustSpeed  \
86232      17.4     29.0       0.0          3.6      11.1
33.0
57576       6.8     14.4       0.8          0.8       8.5
46.0
124071     10.1     15.4       3.2          4.8       8.5
31.0
117955     14.4     33.4       0.0          8.0      11.6
41.0
133468      6.8     14.3       3.2          0.2       7.3
28.0

        WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm ...    NNW
NW   S  \
86232           11.0          19.0         63.0         61.0 ...      0
0   0
```

```
57576            17.0           22.0         80.0        55.0 ...      0
0  1
124071           13.0            9.0         70.0        61.0 ...      0
0  0
117955            9.0           17.0         40.0        23.0 ...      0
0  0
133468           15.0           13.0         92.0        47.0 ...      0
0  0

        SE  SSE  SSW  SW  W  WNW  WSW
86232    0    0    0   0  0    0    0
57576    0    0    0   0  0    0    0
124071   0    1    0   0  0    0    0
117955   0    0    0   1  0    0    0
133468   0    0    0   0  0    0    0

[5 rows x 118 columns]
```

We now have training and testing set ready for model building. Before that, we should map all the feature variables onto the same scale. It is called `feature scaling`. I will do it as follows.

# 11. Feature Scaling

```
X_train.describe()
```

```
              MinTemp         MaxTemp       Rainfall    Evaporation  \
count   113754.000000   113754.000000  113754.000000  113754.000000
mean        12.193497       23.237216       0.675080       5.151606
std          6.388279        7.094149       1.183837       2.823707
min         -8.200000       -4.800000       0.000000       0.000000
25%          7.600000       18.000000       0.000000       4.000000
50%         12.000000       22.600000       0.000000       4.800000
75%         16.800000       28.200000       0.600000       5.400000
max         33.900000       48.100000       3.200000      21.800000

              Sunshine  WindGustSpeed  WindSpeed9am  WindSpeed3pm  \
count   113754.000000   113754.000000  113754.000000  113754.000000
mean         8.041154       39.884074      13.978155      18.614756
std          2.769480       13.116959       8.806558       8.685862
min          0.000000        6.000000       0.000000       0.000000
25%          8.200000       31.000000       7.000000      13.000000
50%          8.500000       39.000000      13.000000      19.000000
75%          8.700000       46.000000      19.000000      24.000000
max         14.500000      135.000000      55.000000      57.000000

            Humidity9am     Humidity3pm      ...                   NNW  \
count   113754.000000   113754.000000      ...         113754.000000
mean        68.867486       51.509547      ...              0.054530
std         18.935587       20.530723      ...              0.227061
min          0.000000        0.000000      ...              0.000000
```

```
25%          57.000000         37.000000         ...            0.000000
50%          70.000000         52.000000         ...            0.000000
75%          83.000000         65.000000         ...            0.000000
max         100.000000        100.000000         ...            1.000000

                     NW                 S                SE               SSE  \
count     113754.000000     113754.000000     113754.000000     113754.000000
mean           0.060288          0.067259          0.101605          0.064059
std            0.238021          0.250471          0.302130          0.244860
min            0.000000          0.000000          0.000000          0.000000
25%            0.000000          0.000000          0.000000          0.000000
50%            0.000000          0.000000          0.000000          0.000000
75%            0.000000          0.000000          0.000000          0.000000
max            1.000000          1.000000          1.000000          1.000000

                    SSW                SW                 W               WNW  \
count     113754.000000     113754.000000     113754.000000     113754.000000
mean           0.056402          0.064464          0.069334          0.060798
std            0.230698          0.245578          0.254022          0.238960
min            0.000000          0.000000          0.000000          0.000000
25%            0.000000          0.000000          0.000000          0.000000
50%            0.000000          0.000000          0.000000          0.000000
75%            0.000000          0.000000          0.000000          0.000000
max            1.000000          1.000000          1.000000          1.000000

                    WSW
count     113754.000000
mean           0.065483
std            0.247378
min            0.000000
25%            0.000000
50%            0.000000
75%            0.000000
max            1.000000

[8 rows x 118 columns]

cols = X_train.columns

from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()

X_train = scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

X_train = pd.DataFrame(X_train, columns=[cols])

X_test = pd.DataFrame(X_test, columns=[cols])
```

```
X_train.describe()

             MinTemp          MaxTemp         Rainfall      Evaporation  \
count  113754.000000   113754.000000   113754.000000   113754.000000
mean        0.484406        0.530004        0.210962        0.236312
std         0.151741        0.134105        0.369949        0.129528
min         0.000000        0.000000        0.000000        0.000000
25%         0.375297        0.431002        0.000000        0.183486
50%         0.479810        0.517958        0.000000        0.220183
75%         0.593824        0.623819        0.187500        0.247706
max         1.000000        1.000000        1.000000        1.000000

             Sunshine    WindGustSpeed    WindSpeed9am     WindSpeed3pm  \
count  113754.000000   113754.000000   113754.000000   113754.000000
mean        0.554562        0.262667        0.254148        0.326575
std         0.190999        0.101682        0.160119        0.152384
min         0.000000        0.000000        0.000000        0.000000
25%         0.565517        0.193798        0.127273        0.228070
50%         0.586207        0.255814        0.236364        0.333333
75%         0.600000        0.310078        0.345455        0.421053
max         1.000000        1.000000        1.000000        1.000000

           Humidity9am     Humidity3pm           ...                 NNW  \
count  113754.000000   113754.000000           ...       113754.000000
mean        0.688675        0.515095           ...            0.054530
std         0.189356        0.205307           ...            0.227061
min         0.000000        0.000000           ...            0.000000
25%         0.570000        0.370000           ...            0.000000
50%         0.700000        0.520000           ...            0.000000
75%         0.830000        0.650000           ...            0.000000
max         1.000000        1.000000           ...            1.000000

                   NW               S              SE             SSE  \
count  113754.000000   113754.000000   113754.000000   113754.000000
mean        0.060288        0.067259        0.101605        0.064059
std         0.238021        0.250471        0.302130        0.244860
min         0.000000        0.000000        0.000000        0.000000
25%         0.000000        0.000000        0.000000        0.000000
50%         0.000000        0.000000        0.000000        0.000000
75%         0.000000        0.000000        0.000000        0.000000
max         1.000000        1.000000        1.000000        1.000000

                  SSW              SW               W             WNW  \
count  113754.000000   113754.000000   113754.000000   113754.000000
mean        0.056402        0.064464        0.069334        0.060798
std         0.230698        0.245578        0.254022        0.238960
min         0.000000        0.000000        0.000000        0.000000
25%         0.000000        0.000000        0.000000        0.000000
50%         0.000000        0.000000        0.000000        0.000000
75%         0.000000        0.000000        0.000000        0.000000
```

```
max          1.000000         1.000000         1.000000         1.000000

                   WSW
count   113754.000000
mean         0.065483
std          0.247378
min          0.000000
25%          0.000000
50%          0.000000
75%          0.000000
max          1.000000

[8 rows x 118 columns]
```

We now have `X_train` dataset ready to be fed into the Logistic Regression classifier. I will do it as follows.

# 12. Model training

```python
# train a logistic regression model on the training set
from sklearn.linear_model import LogisticRegression


# instantiate the model
logreg = LogisticRegression(solver='liblinear', random_state=0)


# fit the model
logreg.fit(X_train, y_train)

LogisticRegression(C=1.0, class_weight=None, dual=False,
fit_intercept=True,
        intercept_scaling=1, max_iter=100, multi_class='warn',
        n_jobs=None, penalty='l2', random_state=0,
solver='liblinear',
        tol=0.0001, verbose=0, warm_start=False)
```

# 13. Predict results

```python
y_pred_test = logreg.predict(X_test)

y_pred_test

array(['No', 'No', 'No', ..., 'No', 'No', 'Yes'], dtype=object)
```

## predict_proba method

**predict_proba** method gives the probabilities for the target variable(0 and 1) in this case, in array form.

`0 is for probability of no rain` and `1 is for probability of rain.`

```python
# probability of getting output as 0 - no rain

logreg.predict_proba(X_test)[:,0]

array([0.91387232, 0.83563172, 0.82035588, ..., 0.97674036,
0.7985333 ,
       0.3073458 ])

# probability of getting output as 1 - rain

logreg.predict_proba(X_test)[:,1]

array([0.08612768, 0.16436828, 0.17964412, ..., 0.02325964,
0.2014667 ,
       0.6926542 ])
```

## 14. Check accuracy score

```python
from sklearn.metrics import accuracy_score

print('Model accuracy score: {0:0.4f}'. format(accuracy_score(y_test,
y_pred_test)))

Model accuracy score: 0.8501
```

Here, **y_test** are the true class labels and **y_pred_test** are the predicted class labels in the test-set.

## Compare the train-set and test-set accuracy

Now, I will compare the train-set and test-set accuracy to check for overfitting.

```python
y_pred_train = logreg.predict(X_train)

y_pred_train

array(['No', 'No', 'No', ..., 'No', 'No', 'No'], dtype=object)

print('Training-set accuracy score: {0:0.4f}'.
format(accuracy_score(y_train, y_pred_train)))

Training-set accuracy score: 0.8476
```

## Check for overfitting and underfitting

```python
# print the scores on training and test set

print('Training set score: {:.4f}'.format(logreg.score(X_train,
y_train)))
```

```
print('Test set score: {:.4f}'.format(logreg.score(X_test, y_test)))

Training set score: 0.8476
Test set score: 0.8501
```

The training-set accuracy score is 0.8476 while the test-set accuracy to be 0.8501. These two values are quite comparable. So, there is no question of overfitting.

In Logistic Regression, we use default value of C = 1. It provides good performance with approximately 85% accuracy on both the training and the test set. But the model performance on both the training and test set are very comparable. It is likely the case of underfitting.

I will increase C and fit a more flexible model.

```
# fit the Logsitic Regression model with C=100

# instantiate the model
logreg100 = LogisticRegression(C=100, solver='liblinear',
random_state=0)


# fit the model
logreg100.fit(X_train, y_train)

LogisticRegression(C=100, class_weight=None, dual=False,
fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l2', random_state=0,
solver='liblinear',
          tol=0.0001, verbose=0, warm_start=False)

# print the scores on training and test set

print('Training set score: {:.4f}'.format(logreg100.score(X_train,
y_train)))

print('Test set score: {:.4f}'.format(logreg100.score(X_test,
y_test)))

Training set score: 0.8478
Test set score: 0.8505
```

We can see that, C=100 results in higher test set accuracy and also a slightly increased training set accuracy. So, we can conclude that a more complex model should perform better.

Now, I will investigate, what happens if we use more regularized model than the default value of C=1, by setting C=0.01.

```
# fit the Logsitic Regression model with C=001
```

```
# instantiate the model
logreg001 = LogisticRegression(C=0.01, solver='liblinear',
random_state=0)


# fit the model
logreg001.fit(X_train, y_train)

LogisticRegression(C=0.01, class_weight=None, dual=False,
fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l2', random_state=0,
solver='liblinear',
          tol=0.0001, verbose=0, warm_start=False)
# print the scores on training and test set

print('Training set score: {:.4f}'.format(logreg001.score(X_train,
y_train)))

print('Test set score: {:.4f}'.format(logreg001.score(X_test,
y_test)))

Training set score: 0.8409
Test set score: 0.8448
```

So, if we use more regularized model by setting C=0.01, then both the training and test set accuracy decrease relatiev to the default parameters.

## Compare model accuracy with null accuracy

So, the model accuracy is 0.8501. But, we cannot say that our model is very good based on the above accuracy. We must compare it with the **null accuracy**. Null accuracy is the accuracy that could be achieved by always predicting the most frequent class.

So, we should first check the class distribution in the test set.

```
# check class distribution in test set

y_test.value_counts()

No      22067
Yes      6372
Name: RainTomorrow, dtype: int64
```

We can see that the occurences of most frequent class is 22067. So, we can calculate null accuracy by dividing 22067 by total number of occurences.

```
# check null accuracy score
```

```
null_accuracy = (22067/(22067+6372))

print('Null accuracy score: {0:0.4f}'. format(null_accuracy))

Null accuracy score: 0.7759
```

We can see that our model accuracy score is 0.8501 but null accuracy score is 0.7759. So, we can conclude that our Logistic Regression model is doing a very good job in predicting the class labels.

Now, based on the above analysis we can conclude that our classification model accuracy is very good. Our model is doing a very good job in terms of predicting the class labels.

But, it does not give the underlying distribution of values. Also, it does not tell anything about the type of errors our classifer is making.

We have another tool called `Confusion matrix` that comes to our rescue.

# 15. Confusion matrix

A confusion matrix is a tool for summarizing the performance of a classification algorithm. A confusion matrix will give us a clear picture of classification model performance and the types of errors produced by the model. It gives us a summary of correct and incorrect predictions broken down by each category. The summary is represented in a tabular form.

Four types of outcomes are possible while evaluating a classification model performance. These four outcomes are described below:-

**True Positives (TP)** – True Positives occur when we predict an observation belongs to a certain class and the observation actually belongs to that class.

**True Negatives (TN)** – True Negatives occur when we predict an observation does not belong to a certain class and the observation actually does not belong to that class.

**False Positives (FP)** – False Positives occur when we predict an observation belongs to a certain class but the observation actually does not belong to that class. This type of error is called **Type I error.**

**False Negatives (FN)** – False Negatives occur when we predict an observation does not belong to a certain class but the observation actually belongs to that class. This is a very serious error and it is called **Type II error.**

These four outcomes are summarized in a confusion matrix given below.

```
# Print the Confusion Matrix and slice it into four pieces

from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred_test)

print('Confusion matrix\n\n', cm)
```

```
print('\nTrue Positives(TP) = ', cm[0,0])

print('\nTrue Negatives(TN) = ', cm[1,1])

print('\nFalse Positives(FP) = ', cm[0,1])

print('\nFalse Negatives(FN) = ', cm[1,0])
```
Confusion matrix

 [[20892  1175]
 [ 3087  3285]]

True Positives(TP) =  20892

True Negatives(TN) =  3285

False Positives(FP) =  1175

False Negatives(FN) =  3087

The confusion matrix shows `20892 + 3285 = 24177 correct predictions` and `3087 + 1175 = 4262 incorrect predictions`.

In this case, we have

- `True Positives` (Actual Positive:1 and Predict Positive:1) - 20892

- `True Negatives` (Actual Negative:0 and Predict Negative:0) - 3285

- `False Positives` (Actual Negative:0 but Predict Positive:1) - 1175 `(Type I error)`

- `False Negatives` (Actual Positive:1 but Predict Negative:0) - 3087 `(Type II error)`

```
# visualize confusion matrix with seaborn heatmap

cm_matrix = pd.DataFrame(data=cm, columns=['Actual Positive:1', 'Actual Negative:0'],
                                    index=['Predict Positive:1', 'Predict Negative:0'])

sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='YlGnBu')
```
<matplotlib.axes._subplots.AxesSubplot at 0xacc3104f60>

# 16. Classification metrices

## Classification Report

**Classification report** is another way to evaluate the classification model performance. It displays the **precision**, **recall**, **f1** and **support** scores for the model. I have described these terms in later.

We can print a classification report as follows:-

```
from sklearn.metrics import classification_report

print(classification_report(y_test, y_pred_test))
              precision    recall  f1-score   support

          No       0.87      0.95      0.91     22067
         Yes       0.74      0.52      0.61      6372

   micro avg       0.85      0.85      0.85     28439
   macro avg       0.80      0.73      0.76     28439
weighted avg       0.84      0.85      0.84     28439
```

## Classification accuracy

```
TP = cm[0,0]
TN = cm[1,1]
```

```
FP = cm[0,1]
FN = cm[1,0]

# print classification accuracy

classification_accuracy = (TP + TN) / float(TP + TN + FP + FN)

print('Classification accuracy :
{0:0.4f}'.format(classification_accuracy))

Classification accuracy : 0.8501
```

## Classification error

```
# print classification error

classification_error = (FP + FN) / float(TP + TN + FP + FN)

print('Classification error : {0:0.4f}'.format(classification_error))

Classification error : 0.1499
```

## Precision

**Precision** can be defined as the percentage of correctly predicted positive outcomes out of all the predicted positive outcomes. It can be given as the ratio of true positives (TP) to the sum of true and false positives (TP + FP).

So, **Precision** identifies the proportion of correctly predicted positive outcome. It is more concerned with the positive class than the negative class.

Mathematically, precision can be defined as the ratio of `TP to (TP + FP)`.

```
# print precision score

precision = TP / float(TP + FP)


print('Precision : {0:0.4f}'.format(precision))

Precision : 0.9468
```

## Recall

Recall can be defined as the percentage of correctly predicted positive outcomes out of all the actual positive outcomes. It can be given as the ratio of true positives (TP) to the sum of true positives and false negatives (TP + FN). **Recall** is also called **Sensitivity**.

**Recall** identifies the proportion of correctly predicted actual positives.

Mathematically, recall can be given as the ratio of `TP to (TP + FN).`

```
recall = TP / float(TP + FN)

print('Recall or Sensitivity : {0:0.4f}'.format(recall))

Recall or Sensitivity : 0.8713
```

## True Positive Rate

**True Positive Rate** is synonymous with **Recall**.

```
true_positive_rate = TP / float(TP + FN)

print('True Positive Rate : {0:0.4f}'.format(true_positive_rate))

True Positive Rate : 0.8713
```

## False Positive Rate

```
false_positive_rate = FP / float(FP + TN)

print('False Positive Rate : {0:0.4f}'.format(false_positive_rate))

False Positive Rate : 0.2635
```

## Specificity

```
specificity = TN / (TN + FP)

print('Specificity : {0:0.4f}'.format(specificity))

Specificity : 0.7365
```

## f1-score

**f1-score** is the weighted harmonic mean of precision and recall. The best possible **f1-score** would be 1.0 and the worst would be 0.0. **f1-score** is the harmonic mean of precision and recall. So, **f1-score** is always lower than accuracy measures as they embed precision and recall into their computation. The weighted average of `f1-score` should be used to compare classifier models, not global accuracy.

## Support

**Support** is the actual number of occurrences of the class in our dataset.

# 17. Adjusting the threshold level

```
# print the first 10 predicted probabilities of two classes- 0 and 1
```

```
y_pred_prob = logreg.predict_proba(X_test)[0:10]

y_pred_prob

array([[0.91387232, 0.08612768],
       [0.83563172, 0.16436828],
       [0.82035588, 0.17964412],
       [0.99025882, 0.00974118],
       [0.95726809, 0.04273191],
       [0.97994232, 0.02005768],
       [0.17838588, 0.82161412],
       [0.23482434, 0.76517566],
       [0.90050811, 0.09949189],
       [0.85480088, 0.14519912]])
```

## Observations

- In each row, the numbers sum to 1.

- There are 2 columns which correspond to 2 classes - 0 and 1.

    – Class 0 - predicted probability that there is no rain tomorrow.

    – Class 1 - predicted probability that there is rain tomorrow.

- Importance of predicted probabilities

    – We can rank the observations by probability of rain or no rain.
- predict_proba process

    – Predicts the probabilities

    – Choose the class with the highest probability

- Classification threshold level

    – There is a classification threshold level of 0.5.

    – Class 1 - probability of rain is predicted if probability > 0.5.

    – Class 0 - probability of no rain is predicted if probability < 0.5.

```
# store the probabilities in dataframe

y_pred_prob_df = pd.DataFrame(data=y_pred_prob, columns=['Prob of - No
rain tomorrow (0)', 'Prob of - Rain tomorrow (1)'])

y_pred_prob_df

   Prob of - No rain tomorrow (0)  Prob of - Rain tomorrow (1)
0                        0.913872                     0.086128
```

```
1                         0.835632                    0.164368
2                         0.820356                    0.179644
3                         0.990259                    0.009741
4                         0.957268                    0.042732
5                         0.979942                    0.020058
6                         0.178386                    0.821614
7                         0.234824                    0.765176
8                         0.900508                    0.099492
9                         0.854801                    0.145199
```

```python
# print the first 10 predicted probabilities for class 1 - Probability
of rain

logreg.predict_proba(X_test)[0:10, 1]
```

```
array([0.08612768, 0.16436828, 0.17964412, 0.00974118, 0.04273191,
       0.02005768, 0.82161412, 0.76517566, 0.09949189, 0.14519912])
```

```python
# store the predicted probabilities for class 1 - Probability of rain

y_pred1 = logreg.predict_proba(X_test)[:, 1]
```

```python
# plot histogram of predicted probabilities


# adjust the font size
plt.rcParams['font.size'] = 12


# plot histogram with 10 bins
plt.hist(y_pred1, bins = 10)


# set the title of predicted probabilities
plt.title('Histogram of predicted probabilities of rain')


# set the x-axis limit
plt.xlim(0,1)


# set the title
plt.xlabel('Predicted probabilities of rain')
plt.ylabel('Frequency')
```

```
Text(0,0.5,'Frequency')
```

## Histogram of predicted probabilities of rain



## Observations

- We can see that the above histogram is highly positive skewed.

- The first column tell us that there are approximately 15000 observations with probability between 0.0 and 0.1.

- There are small number of observations with probability > 0.5.

- So, these small number of observations predict that there will be rain tomorrow.

- Majority of observations predict that there will be no rain tomorrow.

## Lower the threshold

```python
from sklearn.preprocessing import binarize

for i in range(1,5):

    cm1=0

    y_pred1 = logreg.predict_proba(X_test)[:,1]

    y_pred1 = y_pred1.reshape(-1,1)

    y_pred2 = binarize(y_pred1, i/10)

    y_pred2 = np.where(y_pred2 == 1, 'Yes', 'No')
```

```python
    cm1 = confusion_matrix(y_test, y_pred2)

    print ('With',i/10,'threshold the Confusion Matrix is ','\n\
n',cm1,'\n\n',

            'with',cm1[0,0]+cm1[1,1],'correct predictions, ', '\n\n',

            cm1[0,1],'Type I errors( False Positives), ','\n\n',

            cm1[1,0],'Type II errors( False Negatives), ','\n\n',

            'Accuracy score: ', (accuracy_score(y_test, y_pred2)), '\n\
n',

            'Sensitivity: ',cm1[1,1]/(float(cm1[1,1]+cm1[1,0])), '\n\
n',

            'Specificity: ',cm1[0,0]/(float(cm1[0,0]+cm1[0,1])),'\n\n',

            '===============================================', '\
n\n')
```

With 0.1 threshold the Confusion Matrix is

 [[12726  9341]
 [  547  5825]]

 with 18551 correct predictions,

 9341 Type I errors( False Positives),

 547 Type II errors( False Negatives),

 Accuracy score:  0.6523084496641935

 Sensitivity:  0.9141556811048337

 Specificity:  0.5766982371867494

 ======================================================

With 0.2 threshold the Confusion Matrix is

 [[17067  5000]
 [ 1233  5139]]

 with 22206 correct predictions,

 5000 Type I errors( False Positives),

1233 Type II errors( False Negatives),

Accuracy score:  0.7808291430781673

Sensitivity:  0.806497175141243

Specificity:  0.7734173199800607

====================================================

With 0.3 threshold the Confusion Matrix is

 [[19080  2987]
 [ 1873  4499]]

 with 23579 correct predictions,

 2987 Type I errors( False Positives),

 1873 Type II errors( False Negatives),

 Accuracy score:  0.8291079151868912

 Sensitivity:  0.7060577526679221

 Specificity:  0.8646395069560883

====================================================

With 0.4 threshold the Confusion Matrix is

 [[20191  1876]
 [ 2517  3855]]

 with 24046 correct predictions,

 1876 Type I errors( False Positives),

 2517 Type II errors( False Negatives),

 Accuracy score:  0.845529027040332

 Sensitivity:  0.6049905838041432

 Specificity:  0.9149861784565188

====================================================

## Comments

- In binary problems, the threshold of 0.5 is used by default to convert predicted probabilities into class predictions.

- Threshold can be adjusted to increase sensitivity or specificity.

- Sensitivity and specificity have an inverse relationship. Increasing one would always decrease the other and vice versa.

- We can see that increasing the threshold level results in increased accuracy.

- Adjusting the threshold level should be one of the last step you do in the model-building process.

# 18. ROC - AUC

## ROC Curve

Another tool to measure the classification model performance visually is **ROC Curve**. ROC Curve stands for **Receiver Operating Characteristic Curve**. An **ROC Curve** is a plot which shows the performance of a classification model at various classification threshold levels.

The **ROC Curve** plots the **True Positive Rate (TPR)** against the **False Positive Rate (FPR)** at various threshold levels.

**True Positive Rate (TPR)** is also called **Recall**. It is defined as the ratio of `TP to (TP + FN).`

**False Positive Rate (FPR)** is defined as the ratio of `FP to (FP + TN).`

In the ROC Curve, we will focus on the TPR (True Positive Rate) and FPR (False Positive Rate) of a single point. This will give us the general performance of the ROC curve which consists of the TPR and FPR at various threshold levels. So, an ROC Curve plots TPR vs FPR at different classification threshold levels. If we lower the threshold levels, it may result in more items being classified as positve. It will increase both True Positives (TP) and False Positives (FP).

```
# plot ROC Curve

from sklearn.metrics import roc_curve

fpr, tpr, thresholds = roc_curve(y_test, y_pred1, pos_label = 'Yes')

plt.figure(figsize=(6,4))

plt.plot(fpr, tpr, linewidth=2)

plt.plot([0,1], [0,1], 'k--' )

plt.rcParams['font.size'] = 12

plt.title('ROC curve for RainTomorrow classifier')
```
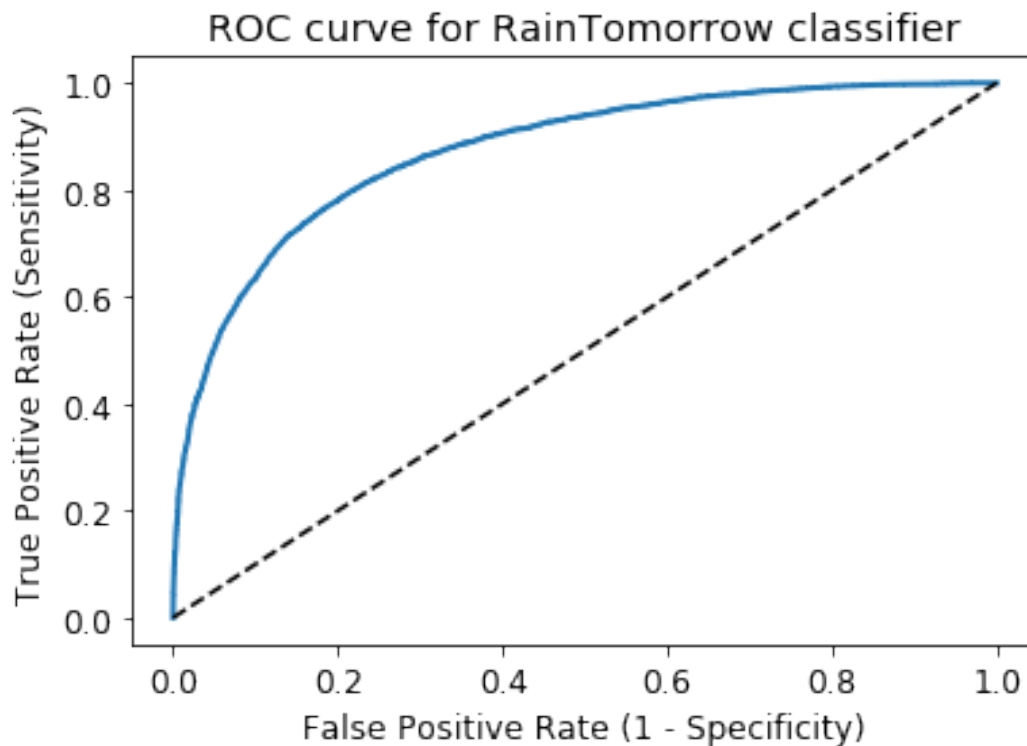
```
plt.xlabel('False Positive Rate (1 - Specificity)')

plt.ylabel('True Positive Rate (Sensitivity)')

plt.show()
```



ROC curve help us to choose a threshold level that balances sensitivity and specificity for a particular context.

## ROC AUC

**ROC AUC** stands for **Receiver Operating Characteristic - Area Under Curve**. It is a technique to compare classifier performance. In this technique, we measure the `area under the curve (AUC)`. A perfect classifier will have a ROC AUC equal to 1, whereas a purely random classifier will have a ROC AUC equal to 0.5.

So, **ROC AUC** is the percentage of the ROC plot that is underneath the curve.

```
# compute ROC AUC

from sklearn.metrics import roc_auc_score

ROC_AUC = roc_auc_score(y_test, y_pred1)

print('ROC AUC : {:.4f}'.format(ROC_AUC))
```

```
ROC AUC : 0.8729
```

## Comments

- ROC AUC is a single number summary of classifier performance. The higher the value, the better the classifier.

- ROC AUC of our model approaches towards 1. So, we can conclude that our classifier does a good job in predicting whether it will rain tomorrow or not.

```python
# calculate cross-validated ROC AUC

from sklearn.model_selection import cross_val_score

Cross_validated_ROC_AUC = cross_val_score(logreg, X_train, y_train,
cv=5, scoring='roc_auc').mean()

print('Cross validated ROC AUC :
{:.4f}'.format(Cross_validated_ROC_AUC))

Cross validated ROC AUC : 0.8695
```

# Model evaluation and improvement

In this section, I will employ several techniques to improve the model performance. I will discuss 3 techniques which are used in practice for performance improvement. These are `recursive feature elimination`, `k-fold cross validation` and `hyperparameter optimization using GridSearchCV`.

# 19. Recursive Feature Elimination with Cross Validation

`Recursive feature elimination (RFE)` is a feature selection technique that helps us to select best features from the given number of features. At first, the model is built on all the given features. Then, it removes the least useful predictor and build the model again. This process is repeated until all the unimportant features are removed from the model.

`Recursive Feature Elimination with Cross-Validated (RFECV) feature selection` technique selects the best subset of features for the estimator by removing 0 to N features iteratively using recursive feature elimination. Then it selects the best subset based on the accuracy or cross-validation score or roc-auc of the model. Recursive feature elimination technique eliminates n features from a model by fitting the model multiple times and at each step, removing the weakest features.

I will use this technique to select best features from this model.

```python
from sklearn.feature_selection import RFECV

rfecv = RFECV(estimator=logreg, step=1, cv=5, scoring='accuracy')
```

```
rfecv = rfecv.fit(X_train, y_train)

print("Optimal number of features : %d" % rfecv.n_features_)

Optimal number of features : 112

# transform the training data

X_train_rfecv = rfecv.transform(X_train)


# train classifier

logreg.fit(X_train_rfecv, y_train)

LogisticRegression(C=1.0, class_weight=None, dual=False,
fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l2', random_state=0,
solver='liblinear',
          tol=0.0001, verbose=0, warm_start=False)

# test classifier on test data

X_test_rfecv = rfecv.transform(X_test)

y_pred_rfecv = logreg.predict(X_test_rfecv)

# print mean accuracy on transformed test data and labels

print ("Classifier score:
{:.4f}".format(logreg.score(X_test_rfecv,y_test)))

Classifier score: 0.8500
```

Our original model accuracy score is 0.8501 whereas accuracy score after RFECV is 0.8500. So, we can obtain approximately similar accuracy but with reduced or optimal set of features.

## Confusion-matrix revisited

I will again plot the confusion-matrix for this model to get an idea of errors our model is making.

```
from sklearn.metrics import confusion_matrix

cm1 = confusion_matrix(y_test, y_pred_rfecv)

print('Confusion matrix\n\n', cm1)

print('\nTrue Positives(TP1) = ', cm1[0,0])

print('\nTrue Negatives(TN1) = ', cm1[1,1])
```

```
print('\nFalse Positives(FP1) = ', cm1[0,1])

print('\nFalse Negatives(FN1) = ', cm1[1,0])

Confusion matrix

 [[20893  1174]
 [ 3091  3281]]

True Positives(TP1) =  20893

True Negatives(TN1) =  3281

False Positives(FP1) =  1174

False Negatives(FN1) =  3091
```

We can see that in the original model, we have FP = 1175 whereas FP1 = 1174. So, we get approximately same number of false positives. Also, FN = 3087 whereas FN1 = 3091. So, we get slightly higher false negatives.

# 20. k-Fold Cross Validation

```
# Applying 10-Fold Cross Validation

from sklearn.model_selection import cross_val_score

scores = cross_val_score(logreg, X_train, y_train, cv = 5,
scoring='accuracy')

print('Cross-validation scores:{}'.format(scores))

Cross-validation scores:[0.84690783 0.84624852 0.84633642 0.84958903
0.84773626]
```

We can summarize the cross-validation accuracy by calculating its mean.

```
# compute Average cross-validation score

print('Average cross-validation score: {:.4f}'.format(scores.mean()))

Average cross-validation score: 0.8474
```

Our, original model score is found to be 0.8476. The average cross-validation score is 0.8474. So, we can conclude that cross-validation does not result in performance improvement.

# 21. Hyperparameter Optimization using GridSearch CV

```python
from sklearn.model_selection import GridSearchCV


parameters = [{'penalty':['l1','l2']},
              {'C':[1, 10, 100, 1000]}]


grid_search = GridSearchCV(estimator = logreg,
                           param_grid = parameters,
                           scoring = 'accuracy',
                           cv = 5,
                           verbose=0)


grid_search.fit(X_train, y_train)

GridSearchCV(cv=5, error_score='raise-deprecating',
        estimator=LogisticRegression(C=1.0, class_weight=None,
dual=False, fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l2', random_state=0,
solver='liblinear',
          tol=0.0001, verbose=0, warm_start=False),
        fit_params=None, iid='warn', n_jobs=None,
        param_grid=[{'penalty': ['l1', 'l2']}, {'C': [1, 10, 100,
1000]}],
        pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
        scoring='accuracy', verbose=0)
# examine the best model

# best score achieved during the GridSearchCV
print('GridSearch CV best score : {:.4f}\n\
n'.format(grid_search.best_score_))

# print parameters that give the best results
print('Parameters that give the best results :','\n\n',
(grid_search.best_params_))

# print estimator that was chosen by the GridSearch
print('\n\nEstimator that was chosen by the search :','\n\n',
(grid_search.best_estimator_))

GridSearch CV best score : 0.8474


Parameters that give the best results :
```

```
 {'penalty': 'l1'}


Estimator that was chosen by the search :

 LogisticRegression(C=1.0, class_weight=None, dual=False,
fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l1', random_state=0,
solver='liblinear',
          tol=0.0001, verbose=0, warm_start=False)
```

```python
# calculate GridSearch CV score on test set

print('GridSearch CV score on test set:
{0:0.4f}'.format(grid_search.score(X_test, y_test)))
```

```
GridSearch CV score on test set: 0.8507
```

## Comments

- Our original model test accuracy is 0.8501 while GridSearch CV accuracy is 0.8507.

- We can see that GridSearch CV improve the performance for this particular model.

# 22. Results and Conclusion

1. The logistic regression model accuracy score is 0.8501. So, the model does a very good job in predicting whether or not it will rain tomorrow in Australia.

2. Small number of observations predict that there will be rain tomorrow. Majority of observations predict that there will be no rain tomorrow.

3. The model shows no signs of overfitting.

4. Increasing the value of C results in higher test set accuracy and also a slightly increased training set accuracy. So, we can conclude that a more complex model should perform better.

5. Increasing the threshold level results in increased accuracy.

6. ROC AUC of our model approaches towards 1. So, we can conclude that our classifier does a good job in predicting whether it will rain tomorrow or not.

7. Our original model accuracy score is 0.8501 whereas accuracy score after RFECV is 0.8500. So, we can obtain approximately similar accuracy but with reduced set of features.

8. In the original model, we have FP = 1175 whereas FP1 = 1174. So, we get approximately same number of false positives. Also, FN = 3087 whereas FN1 = 3091. So, we get slighly higher false negatives.

9. Our, original model score is found to be 0.8476. The average cross-validation score is 0.8474. So, we can conclude that cross-validation does not result in performance improvement.

10. Our original model test accuracy is 0.8501 while GridSearch CV accuracy is 0.8507. We can see that GridSearch CV improve the performance for this particular model.