

```
[2]: import pandas as pd
import matplotlib.pyplot as plt
import pickle
import matplotlib inline
import warnings
warnings.filterwarnings('ignore')
from pandas.plotting import scatter_matrix
from sklearn.preprocessing import LabelEncoder
import seaborn as sns
import numpy as np

In [2]: data=pd.read_csv("weatherAUS.csv")
data.head()
```

	Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindSpeed9am	Humidity9am	Humidity3pm	Pressure9am	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow	
0	2008-12-03	Albury	13.4	22.9	0.6		NaN	NaN	W	44.0	W	...	71.0	22.0	1007.7	1007.1	8.0	NaN	16.9	21.8	No	
1	2008-12-02	Albury	7.4	25.1	0.0		NaN	NaN	WNN	44.0	NNW	...	44.0	25.0	1010.6	1007.8	NaN	NaN	17.2	24.3	No	
2	2008-12-04	Albury	12.9	25.7	0.0		NaN	NaN	WSW	46.0	W	...	38.0	30.0	1007.6	1008.7	NaN	2.0	21.0	23.2	No	
3	2008-12-04	Albury	9.2	28.0	0.0		NaN	NaN	NE	24.0	SE	...	45.0	16.0	1017.6	1012.8	NaN	NaN	18.1	26.5	No	
4	2008-12-05	Albury	17.5	32.3	1.0		NaN	NaN	W	41.0	ENE	...	82.0	33.0	1010.8	1006.0	7.0	8.0	17.8	29.7	No	

5 rows × 23 columns

```
In [4]: data.shape
Out[4]: (8425, 23)

In [5]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8425 entries, 0 to 8424
Data columns (total 23 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Date                   8425 non-null   object
1   Location               8425 non-null   object
2   MinTemp                8350 non-null   float64
3   MaxTemp                8365 non-null   float64
4   Rainfall              8185 non-null   float64
5   Evaporation            4913 non-null   float64
6   Sunshine               4433 non-null   float64
7   WindGustDir            7434 non-null   object
8   WindGustSpeed          7434 non-null   float64
9   WindDir9am             7596 non-null   object
10  WindDir3pm             8117 non-null   object
11  WindSpeed9am           8349 non-null   float64
12  WindSpeed3pm           8318 non-null   float64
13  Humidity9am            8366 non-null   float64
14  Humidity3pm            8323 non-null   float64
15  Pressure9am            7116 non-null   float64
16  Pressure3pm            7113 non-null   float64
17  Cloud9am                6804 non-null   float64
18  Cloud3pm               5970 non-null   float64
19  Temp9am                8369 non-null   float64
20  Temp3pm                8329 non-null   float64
21  RainToday              8185 non-null   object
22  RainTomorrow           8186 non-null   object
dtypes: float64(16), object(7)
memory usage: 1.9+Mib

In [6]: data.isnull().sum()
Out[6]: Date                0
Location                0
MinTemp                75
MaxTemp                60
Rainfall              240
Evaporation            3512
Sunshine               3994
WindGustDir            991
WindGustSpeed          829
WindDir9am             996
WindDir3pm             998
WindSpeed9am           76
WindSpeed3pm           187
Humidity9am            59
Humidity3pm            182
Pressure9am            1309
Pressure3pm            1312
Cloud9am               2421
Cloud3pm               2455
Temp9am                96
Temp3pm                240
RainToday              239
RainTomorrow           239
dtype: int64

In [7]: data.describe()
Out[7]:
      MinTemp  MaxTemp  Rainfall  Evaporation  Sunshine  WindGustSpeed  WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm  Pressure9am  Pressure3pm  Cloud9am  Cloud3pm  Temp9am  Temp3pm
count  8350.000000  8365.000000  8185.000000  4913.000000  4433.000000  7434.000000  8349.000000  8318.000000  8366.000000  8323.000000  7116.000000  7113.000000  6804.000000  5970.000000  8369.000000  8329.000000  8329.000000
mean    13.93305   23.69976    2.80913    5.390995    7.632205    46.174469    13.847646    18.523662    67.822496    51.249790    1017.640233    1015.236075    4.566622    4.503183    17.782015    22.442394
std     6.540356    6.126408    10.49379    5.044484    3.986225    14.665721    10.174579    9.766966    16.832383    18.423774    6.828699    6.766681    2.877658    2.731659    5.627025    5.960300
min     2.000000    8.200000    0.000000    0.000000    0.000000    7.000000    0.000000    0.000000    10.000000    6.000000    989.800000    982.800000    0.000000    0.000000    1.900000    7.300000
25%     8.200000    19.300000    0.000000    2.600000    4.750000    6.000000    11.000000    6.000000    39.000000    39.000000    1013.000000    1010.400000    1.000000    2.000000    13.800000    18.000000
50%    13.000000    23.300000    0.000000    4.600000    8.700000    39.000000    13.000000    19.000000    68.000000    51.000000    1017.700000    1015.300000    5.000000    5.000000    17.800000    21.900000
75%    17.400000    28.000000    1.000000    7.000000    10.700000    50.000000    20.000000    24.000000    80.000000    63.000000    1022.300000    1019.800000    7.000000    7.000000    21.900000    26.400000
max     28.500000    45.000000    371.000000    145.000000    13.900000    107.000000    63.000000    83.000000    100.000000    98.000000    1038.000000    1036.000000    8.000000    8.000000    38.400000    44.300000

In [8]: data=data.drop(["Date","Location","Evaporation","Sunshine","Cloud9am","Cloud3pm"],axis =1)
data.head()
Out[8]:
      MinTemp  MaxTemp  Rainfall  WindGustDir  WindGustSpeed  WindDir9am  WindDir3pm  WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm  Pressure9am  Pressure3pm  Temp9am  Temp3pm  RainToday  RainTomorrow
0      13.4      22.9      0.6           W           44.0           W          WSW           20.0           24.0           71.0           22.0           1007.7           1007.1           16.9           21.8           No           No
1       7.4      25.1      0.0          WNW           44.0          NNW          WSW           4.0           22.0           44.0           25.0           1010.6           1007.8           17.2           24.3           No           No
2      12.9      25.7      0.0          WSW           46.0           W          WSW           19.0           26.0           38.0           30.0           1007.6           1008.7           21.0           23.2           No           No
3       9.2      28.0      0.0           NE           24.0          SE           E           11.0           9.0           45.0           16.0           1017.6           1012.8           18.1           26.5           No           No
4      17.5      32.3      1.0           W           41.0          ENE          NW           7.0           20.0           82.0           33.0           1010.8           1006.0           17.8           29.7           No           No

In [9]: data.columns
Out[9]:
Index(['MinTemp', 'MaxTemp', 'Rainfall', 'WindGustDir', 'WindGustSpeed',
       'WindDir9am', 'WindDir3pm', 'WindSpeed9am', 'WindSpeed3pm',
       'Humidity9am', 'Humidity3pm', 'Pressure9am', 'Pressure3pm', 'Temp9am',
       'Temp3pm', 'RainToday', 'RainTomorrow'],
      dtype='object')

In [10]: le =LabelEncoder()
data["WindGustDir"]=le.fit_transform(data["WindGustDir"])
data["WindDir9am"]=le.fit_transform(data["WindDir9am"])
data["WindDir3pm"]=le.fit_transform(data["WindDir3pm"])
data["RainToday"]=le.fit_transform(data["RainToday"])
data["RainTomorrow"]=le.fit_transform(data["RainTomorrow"])
data.head()
Out[10]:
      MinTemp  MaxTemp  Rainfall  WindGustDir  WindGustSpeed  WindDir9am  WindDir3pm  WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm  Pressure9am  Pressure3pm  Temp9am  Temp3pm  RainToday  RainTomorrow
0      13.4      22.9      0.6           13           44.0           13          14           20.0           24.0           71.0           22.0           1007.7           1007.1           16.9           21.8           0           0
1       7.4      25.1      0.0           14           44.0           6          15          4.0           22.0           44.0           25.0           1010.6           1007.8           17.2           24.3           0           0
2      12.9      25.7      0.0           15           46.0           13          15          19.0           26.0           38.0           30.0           1007.6           1008.7           21.0           23.2           0           0
3       9.2      28.0      0.0           4          24.0           9           0          11.0           9.0           45.0           16.0           1017.6           1012.8           18.1           26.5           0           0
4      17.5      32.3      1.0           13           41.0           1          7           7.0           20.0           82.0           33.0           1010.8           1006.0           17.8           29.7           0           0

In [12]: data=data.dropna(axis =0)

In [13]: data.isnull().sum()
Out[13]:
MinTemp      0
MaxTemp      0
Rainfall      0
WindGustDir   0
WindGustSpeed 0
WindDir9am    0
WindDir3pm    0
WindSpeed9am  0
WindSpeed3pm  0
Humidity9am   0
Humidity3pm   0
Pressure9am   0
Pressure3pm   0
Temp9am       0
Temp3pm       0
RainToday     0
RainTomorrow   0
dtype: int64

In [15]: plt.figure(figsize=(20,20))
plotnumber=1
for column in data:
    if plotnumber<=17:
        ax=plt.subplot(6,3,plotnumber)
        sns.distplot(data[column])
        plt.xlabel(column,fontsize=20)
        plotnumber+=1
    else:
        plotnumber+=1
plt.tight_layout()

Density
0.07
0.06
0.05
0.04
0.03
0.02
0.01
0.00
-5 0 5 10 15 20 25 30
MinTemp

Density
0.07
0.06
0.05
0.04
0.03
0.02
0.01
0.00
10 20 30 40 50
MaxTemp

Density
0.200
0.175
0.150
0.125
0.100
0.075
0.050
0.025
0.000
0 50 100 150 200
Rainfall

Density
0.10
0.08
0.06
0.04
0.02
0.00
-2.5 0.0 2.5 5.0 7.5 10.0 12.5 15.0
WindGustDir

Density
0.05
0.04
0.03
0.02
0.01
0.00
0 20 40 60 80 100 120
WindGustSpeed

Density
0.12
0.10
0.08
0.06
0.04
0.02
0.00
0 5 10 15
WindDir9am

Density
0.10
0.08
0.06
0.04
0.02
0.00
0 5 10 15
WindDir3pm

Density
0.07
0.06
0.05
0.04
0.03
0.02
0.01
0.00
0 10 20 30 40 50 60 70
WindSpeed9am

Density
0.06
0.05
0.04
0.03
0.02
0.01
0.00
0 20 40 60 80 100
WindSpeed3pm

Density
0.06
0.05
0.04
0.03
0.02
0.01
0.00
980 990 1000 1010 1020 1030 1040
Humidity9am

Density
0.07
0.06
0.05
0.04
0.03
0.02
0.01
0.00
0 20 40 60 80 100
Humidity3pm

Density
0.07
0.06
0.05
0.04
0.03
0.02
0.01
0.00
5 10 15 20 30 40
Temp9am

Density
0.07
0.06
0.05
0.04
0.03
0.02
0.01
0.00
10 20 30 40
Temp3pm

Density
35
30
25
20
15
10
5
0
-0.2 0.0 0.2 0.4 0.6 0.8 1.0 1.2
RainToday

Density
17.5
15.0
12.5
10.0
7.5
5.0
2.5
0.0
0.0 0.5 1.0 1.5 2.0
RainTomorrow

In [18]: scatter_matrix(data, figsize=(45,45),diagonal="hist",color="orange")
plt.show()
```

```
In [34]: from sklearn.metrics import classification_report
import statsmodels.formula.api as smf
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, roc_auc_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier

In [20]: x=data.drop(['RainTomorrow'], axis =1)
y=data['RainTomorrow']

In [22]: plt.figure(figsize=(10,10)) , hue = "RainTomorrow" , data =data)
plt.show()
```

```
In [23]: plt.figure(figsize=(20,25))
graph=plt
for column in x:
    if graph==6:
        ax=plt.subplot(6,2,graph)
        sns.boxplot(x=column)
        plt.xlabel(column,fontsize=10)
        graph+=1
    else:
        graph+=1
plt.show()
```

```
In [24]: x_train,x_test,y_train,y_test= train_test_split(x,y, test_size=0.2)

In [25]: lr =LogisticRegression()
lr.fit(x_train,y_train)

Out[25]: LogisticRegression()

In [26]: y_pred =lr.predict(x_test)
y_pred
Out[26]: array([1, 0, 0, ..., 0, 0, 1])

In [29]: print("Accuracy =",accuracy_score(y_test,y_pred))
print("Confusion matrix =",confusion_matrix(y_test,y_pred))
Accuracy = 0.8303571428571429
Confusion matrix = [[946 65 0]
 [14 236 9]
 [19 3 0]]

In [30]: # classifier report
print(classification_report(y_test,y_pred))
precision    recall  f1-score   support

0   0.86      0.86      0.89      1011
1   0.72      0.76      0.74      311
2   0.68      0.68      0.68       22

accuracy          0.83      1344
macro avg        0.72      0.78      0.78      1344
weighted avg     0.81      0.83      0.82      1344

In [38]: # Decision Tree Classifier
dtc =DecisionTreeClassifier()
dtc.fit(x_train,y_train)

In [38]: dtc.predict(x_test)

In [39]: y_predi =dtc.predict(x_test)
y_predi
Out[39]: array([1, 0, 0, ..., 0, 0, 1])

In [40]: print("Accuracy =",accuracy_score(y_test,y_predi))
print("Confusion matrix =",confusion_matrix(y_test,y_predi))
Accuracy = 0.8683035714285714
Confusion matrix = [[916 89 6]
 [ 4 236 1]
 [ 3 4 15]]

In [43]: print(classification_report(y_test,y_predi))
precision    recall  f1-score   support

0   0.92      0.91      0.94      1011
1   0.72      0.76      0.74      311
2   0.68      0.68      0.68       22

accuracy          0.83      1344
macro avg        0.77      0.78      0.78      1344
weighted avg     0.87      0.87      0.87      1344

In [ ]:

In [42]: # Random Forest Classifier
rf =RandomForestClassifier()
rf.fit(x_train,y_train)

Out[42]: RandomForestClassifier()

In [42]: rf.predict(x_test)
y_predict
Out[42]: array([1, 0, 0, ..., 0, 0, 1])

In [43]: print("Accuracy =",accuracy_score(y_test,y_predict))
print("Confusion matrix =",confusion_matrix(y_test,y_predict))
Accuracy = 0.9055059523809523
Confusion matrix = [[978 33 0]
 [ 8 224 6]
 [ 5 2 15]]

In [46]: print(classification_report(y_test,y_predict))
precision    recall  f1-score   support

0   0.91      0.97      0.94      1011
1   0.86      0.72      0.79      311
2   1.00      0.68      0.81       22

accuracy          0.93      1344
macro avg        0.93      0.79      0.90      1344
weighted avg     0.90      0.91      0.91      1344

In [ ]: 
```