

```
import pandas as pd
import matplotlib.pyplot as plt
import pickle
import math
import warnings
warnings.filterwarnings('ignore')
import seaborn as sns
from sklearn.linear_model import LinearRegression
import statsmodels.formula.api as smf
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, roc_auc_score
from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, roc_auc_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier

In [2]: data=pd.read_csv('weatherAUS.csv')

Out[3]: data.head()
   Date Location MinTemp MaxTemp Rainfall Evaporation Sunshine WindGustDir WindGustSpeed WindDir9am ... Humidity5am Humidity3pm Pressure9am Pressure3pm Cloud9am Cloud3pm Temp9am Temp3pm RainToday RainTomorrow
0 2006-12-05 Albany 13.4 22.9 0.6 0.0 NaN NaN W 44.0 W ... 71.0 22.0 1007.7 1007.1 8.0 NaN 16.9 21.8 21.8 No
1 2006-12-05 Albany 17.4 25.1 0.0 0.0 NaN NaN WNW 44.0 NWW ... 44.0 25.0 1007.6 1007.8 NaN NaN 17.2 24.3 No
2 2006-12-05 Albany 12.9 25.7 0.0 WSW 46.0 W ... 38.0 30.0 1007.6 1008.7 NaN NaN 21.0 23.2 No
3 2006-12-05 Albany 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 2006-12-05 Albany 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 28.7 No
5 rows x 23 columns

In [4]: data.shape
Out[4]: (8425, 28)

In [5]: data.isnull().sum()
Out[5]: Date 0
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
Temp9am 0
Temp3pm 0
RainToday 0
RainTomorrow 0
dtype: int64

In [6]: data.describe()
Out[6]:
   MinTemp  MaxTemp  Rainfall  Evaporation  Sunshine  WindGustSpeed  WindGustDir  WindDir9am  WindDir3pm  WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm  Pressure9am  Pressure3pm  Temp9am  Temp3pm  RainToday  RainTomorrow
count  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000  8425.000000
mean    13.383363  23.899776  2.859913  5.389395  7.632205  40.174469  13.847646  18.533662  16.822486  51.249790  1017.640233  1015.236075  4.566622  4.503183  17.702015  22.442934
std     5.405360  6.134020  10.459379  5.044844  3.896325  14.665711  10.174579  9.765896  16.833283  18.423774  6.836869  6.766661  2.877058  2.731459  5.637355  5.880232
min     -2.000000  8.200000  0.000000  0.000000  0.000000  7.000000  0.000000  0.000000  0.000000  0.000000  989.800000  982.900000  0.000000  0.000000  1.900000  7.300000
25%     9.200000  19.300000  0.000000  2.000000  7.500000  30.000000  0.000000  11.000000  16.000000  39.000000  1013.000000  1010.400000  1.000000  1.000000  13.800000  18.000000
50%    13.100000  23.300000  0.000000  4.000000  8.700000  39.000000  0.000000  15.000000  20.000000  45.000000  1017.000000  1015.900000  3.000000  3.000000  17.800000  21.900000
75%    17.400000  26.000000  0.000000  6.000000  10.000000  50.000000  0.000000  20.000000  24.000000  60.000000  1022.000000  1019.900000  7.000000  7.000000  21.900000  26.400000
max    21.000000  45.500000  371.000000  145.000000  13.900000  107.000000  63.000000  83.000000  100.000000  90.000000  1029.000000  1026.000000  8.000000  8.000000  39.400000  44.100000

In [7]: data.info()
Out[7]:
<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             8425 non-null    float64
 3   MaxTemp             8425 non-null    float64
 4   Rainfall            8425 non-null    float64
 5   Evaporation          8425 non-null    float64
 6   Sunshine            8425 non-null    float64
 7   WindGustDir         8425 non-null    object
 8   WindGustSpeed       8425 non-null    float64
 9   WindDir9am          8425 non-null    float64
10  WindDir3pm          8425 non-null    float64
11  WindSpeed9am        8425 non-null    float64
12  WindSpeed3pm        8425 non-null    float64
13  Humidity9am         8425 non-null    float64
14  Humidity3pm         8425 non-null    float64
15  Pressure9am         8425 non-null    float64
16  Pressure3pm         8425 non-null    float64
17  Temp9am             8425 non-null    float64
18  Temp3pm             8425 non-null    float64
19  RainToday           8425 non-null    object
20  RainTomorrow         8425 non-null    object
dtypes: float64(16), object(7)
memory usage: 1.5+ MB

In [8]: data=data.drop(['Date','Location','Evaporation','Sunshine','Cloud9am','Cloud3pm','Temp9am','Temp3pm','RainToday','RainTomorrow'],axis=1)
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          44.0         W          20.0    24.0    71.0    22.0    1007.7    1007.1    16.9    21.8    0         0
1    17.4    25.1    0.0        WNW         44.0         NWW         WSW         4.0    22.0    44.0    25.0    1007.6    1007.8    17.2    24.3    0         0
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    0         0
3     9.2    28.0    0.0         NE         24.0         SE         SE         11.0    16.0    45.0    16.0    1017.6    1012.8    18.1    26.5    0         0
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    28.7    0         0

In [9]: data=data.dropna(axis=6)
Out[9]: data.shape
Out[10]: (8332, 17)

In [11]: data.columns
Out[11]: Index(['MinTemp', 'MaxTemp', 'Rainfall', 'WindGustDir', 'WindGustSpeed', 'WindDir9am', 'WindDir3pm', 'WindSpeed9am', 'WindSpeed3pm', 'Humidity9am', 'Humidity3pm', 'Pressure9am', 'Pressure3pm', 'Temp9am', 'Temp3pm', 'RainToday', 'RainTomorrow'],
      dtype='object')

In [12]: le=LabelEncoder()

In [13]: 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'])

In [14]: data.head()
Out[14]:
   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    17.4    25.1    0.0        WNW         44.0         NWW         WSW         4.0    22.0    44.0    25.0    1007.6    1007.8    17.2    24.3    0         0
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    0         0
3     9.2    28.0    0.0         NE         24.0         SE         SE         11.0    16.0    45.0    16.0    1017.6    1012.8    18.1    26.5    0         0
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    28.7    0         0

In [15]: data.isnull().sum()
Out[15]:
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 [16]: plt.figure(figsize=(20,20),facecolor='yellow')
plt.plotnumber=1
for column in x:
    if plotnumber<=7:
        ax=plt.subplot(8,2,plotnumber)
        sns.distplot(x[column],y)
        plt.xlabel(column,fontsize=20)
        plotnumber+=1
    plt.tight_layout()

In [17]: x=data.drop(['RainTomorrow'],axis=5)
y=data['RainTomorrow']

In [18]: print(x,y)
Out[18]:
   MinTemp  MaxTemp  Rainfall  WindGustDir  WindGustSpeed  WindDir9am  \
0    13.4    22.9    0.6         13          44.0         13
1    17.4    25.1    0.0        WNW         44.0         NWW
2    12.9    25.7    0.0        WSW         46.0         W
3     9.2    28.0    0.0         NE         24.0         SE
4    17.5    32.3    1.0         W          41.0         ENE
... ..
8419    13.9    23.8    0.8         15          50.0         15
8420    2.8    23.4    0.8         0          31.0         0
8421    1.6    20.3    0.8         3          26.0         3
8422    5.4    26.9    0.8         3          37.0         9
8423    1.8    27.0    0.0         0          26.0         0

   WindDir9am  WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm  \
0    14.0    20.0    24.0    71.0    22.0
1    15.0    22.0    25.0    44.0    25.0
2    15.0    26.0    38.0    30.0    30.0
3    11.0    16.0    45.0    16.0    23.0
4    7.0    20.0    82.0    33.0    33.0
... ..
8419    15.0    50.0    50.0    27.0    27.0
8420    3.0    11.0    11.0    51.0    24.0
8421    3.0    9.0    9.0    26.0    23.0
8422    9.0    9.0    9.0    53.0    24.0
8423    7.0    7.0    7.0    24.0    24.0

   Pressure9am  Pressure3pm  Temp9am  Temp3pm  RainToday
0    1007.7    1007.1    16.9    21.8    0
1    1007.6    1007.8    17.2    24.3    0
2    1007.6    1008.7    21.0    23.2    0
3    1017.6    1012.8    18.1    26.5    0
4    1010.8    1006.0    17.8    28.7    0
... ..
8419    1024.7    1021.2    26.9    26.9    0
8420    1024.6    1020.3    26.1    22.4    0
8421    1023.5    1021.1    26.0    24.5    0
8422    1021.0    1016.8    12.5    26.1    0
8423    1021.4    1015.5    15.1    26.0    0

[6332 rows x 16 columns] 0 0
0
1
8419 0
8420 0
8421 0
8422 0
8423 0
Name: RainTomorrow, Length: 6332, dtype: int32

In [19]: plt.figure(figsize=(15,20),facecolor='yellow')
plt.plotnumber=1
for column in x:
    if plotnumber<=8:
        ax=plt.subplot(8,2,plotnumber)
        plt.scatter(x[column],y)
        plt.xlabel(column,fontsize=16)
        plotnumber+=1
    plt.tight_layout()

In [20]: plt.figure(figsize=(15,20),facecolor='yellow')
plt.plotnumber=1
for column in x:
    if plotnumber<=8:
        ax=plt.subplot(8,2,plotnumber)
        sns.boxplot(x[column],y)
        plt.xlabel(column,fontsize=16)
        plotnumber+=1
    plt.show()

In [21]: plt.figure(figsize=(10,10))
sns.scatterplot(x='Humidity9am',y='Temp9am',hue='RainTomorrow',data=data)
plt.show()

In [22]: plt.figure(figsize=(10,10))
sns.scatterplot(x='MinTemp',y='MaxTemp',hue='RainTomorrow',data=data)
plt.show()

In [23]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2)

In [24]: print(x_train,y_train)
Out[24]:
   MinTemp  MaxTemp  Rainfall  WindGustDir  WindGustSpeed  WindDir9am  \
6513    28.1    29.7    0.8         0          26.0         0
7726    12.3    24.8    0.0         12          36.0         1
1451    18.3    27.9    13.8         3          33.0         6
4313    1.4    21.9    0.2         6          44.0         6
6501    16.8    26.9    0.0         0          50.0         8
... ..
4857    13.7    26.8    0.0         8          65.0    13
4823    12.0    22.6    0.0         7          66.0         7
4823    9.5    13.8    0.0         3          46.0         6
5483    19.8    26.8    1.2    15.0    102.0    12
4836    12.8    18.2    0.8         14          76.0    15

   WindDir9am  WindSpeed9am  WindSpeed3pm  Humidity9am  Humidity3pm  \
6513         1         7.8    11.0    65.0    61.0
7726    12         17.8    11.8    36.0    39.0
1451         6         26.0    39.0    45.0    30.0
4313         4         12.0    9.0    67.0    65.0
6501         2         6.0    15.0    66.0    54.9
... ..
4857         9        30.0    26.0    36.0    32.0
4823    14        15.0    20.0    51.0    21.0
5483    11         6.0    26.0    68.0    74.9
4836    13        20.0    16.0    41.0    31.9

   Pressure9am  Pressure3pm  Temp9am  Temp3pm  RainToday
6513    1026.9    1014.3    26.4    29.1    0
7726    1022.0    1023.5    16.6    22.9    0
1451    1014.7    1011.2    23.0    27.4    1
4313    1024.5    1021.6    19.0    22.4    0
6501    1021.6    1019.9    23.9    25.4    0
... ..
4857    1009.7    1007.0    20.1    22.6    0
4823    1017.0    1013.1    20.4    22.9    0
4823    1016.1    1014.0    19.9    11.7    0
5483    1009.6    1003.2    22.9    18.2    1
4836    1021.1    1012.3    13.9    17.3    0

[5665 rows x 16 columns] 6513 0
7726 0
1451 0
4313 0
6501 0
... ..
4857 0
4823 0
4823 0
5483 0
4836 0
Name: RainTomorrow, Length: 5665, dtype: int32

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

In [26]: y_pred=lr.predict(x_test)

In [27]: array([0, 1, 0, ..., 0, 0, 0])

In [28]: # accuracy
accuracy=accuracy_score(y_test,y_pred)
accuracy
Out[28]: 0.8304348929045

In [29]: # confusion matrix
con_mat=confusion_matrix(y_test,y_pred)
con_mat
Out[29]: array([[891,  67],
       [142, 367]], dtype=int64)

In [30]: from sklearn.metrics import classification_report

In [31]: # classifier report
print(classification_report(y_test,y_pred))
Out[31]:
precision    recall  f1-score   support

0.00         0.86     0.93     958
1.00         0.71     0.84     399

accuracy: 0.79
macro avg: 0.79     0.74     0.76    1267
weighted avg: 0.83     0.84     0.83    1267

In [32]: # ROC Curve
fpr,tpr,threshold=roc_curve(y_test,y_pred)

In [33]: print('threshold:',threshold)
print('True Positive Rate:',tpr)
print('False Positive Rate:',fpr)

threshold: [2 0]
True Positive Rate: [0. 0.5404387 1. ]
False Positive Rate: [0. 0.06993737 1. ]

In [34]: plt.plot(fpr,tpr,color='c',label='ROC')
plt.plot([0,1],[0,1],color='m',linestyle='--')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC Curve)')
plt.legend()

In [35]: Receiver Operating Characteristic (ROC Curve)
True Positive Rate
False Positive Rate
ROC

In [36]: # AUC curve
auc_score=roc_auc_score(y_test,y_pred)
print(auc_score)
Out[36]: 0.7529574234885

In [37]: # Random Forest Classifier
rf=RandomForestClassifier()
rf.fit(x_train,y_train)

In [38]: y_pred=rf.predict(x_test)

In [39]: array([0, 1, 0, ..., 0, 0, 0])

In [40]: # accuracy
accuracy=accuracy_score(y_test,y_pred)
accuracy
Out[40]: 0.910023677974791

In [41]: # confusion matrix
con_mat=confusion_matrix(y_test,y_pred)
con_mat
Out[41]: array([[930,  28],
       [ 86, 223]], dtype=int64)

In [42]: # classifier report
print(classification_report(y_test,y_pred))
Out[42]:
precision    recall  f1-score   support

0.00         0.92     0.91     958
1.00         0.74     0.76     399

accuracy: 0.89
macro avg: 0.89     0.85     0.87    1267
weighted avg: 0.91     0.85     0.91    1267

In [43]: # Decision Tree Classifier
dtc=DecisionTreeClassifier()
dtc.fit(x_train,y_train)

In [44]: y_pred=dtc.predict(x_test)

In [45]: array([1, 1, 0, ..., 0, 0, 0])

In [46]: # accuracy
accuracy=accuracy_score(y_test,y_pred)
accuracy
Out[46]: 0.87259974734885

In [47]: # confusion matrix
con_mat=confusion_matrix(y_test,y_pred)
con_mat
Out[47]: array([[875,  23],
       [ 75, 241]], dtype=int64)

In [48]: # classifier report
print(classification_report(y_test,y_pred))
Out[48]:
precision    recall  f1-score   support

0.00         0.92     0.91     958
1.00         0.74     0.76     399

accuracy: 0.89
macro avg: 0.89     0.85     0.87    1267
weighted avg: 0.91     0.85     0.91    1267

In [49]:
In [50]:
In [51]:
In [52]:
```