

Week 3 Homework - Jon Workman

```
In [36]: from datetime import date
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.dummy import DummyClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
import seaborn as sns
from sklearn import svm, datasets, metrics
from sklearn.metrics import classification_report, confusion_matrix, roc_curve, accuracy_score
%matplotlib inline
```

Load dataset

```
In [37]: # Load Seattle weather dataset
df = pd.read_csv('seattleWeather_1948-2017.csv').dropna()
```

Clean data

```
In [38]: # Clean data
df.RAIN = df.RAIN.astype(int)
df['DATE'] = pd.to_datetime(df['DATE'], format='%Y-%m-%d')
```

Create additional feature

```
In [39]: # Average precipitation for last 2 days feature.
# Initially tried 7 days and 5 days. 2 days was the best predictor of raining on a given day.

df['average_PRCP'] = df['PRCP'].rolling(2).mean().round(2)
```

```
In [40]: # Find any null values
df[df.isnull().any(axis=1)]
```

Out[40]:

	DATE	PRCP	TMAX	TMIN	RAIN	average_PRCP
0	1948-01-01	0.47	51	42	1	NaN

```
In [41]: # Remove null values
df.drop(df.index[0], inplace=True)
```

```
In [42]: df.head(5)
```

Out[42]:

	DATE	PRCP	TMAX	TMIN	RAIN	average_PRCP
1	1948-01-02	0.59	45	36	1	0.53
2	1948-01-03	0.42	45	35	1	0.50
3	1948-01-04	0.31	45	34	1	0.37
4	1948-01-05	0.17	45	32	1	0.24
5	1948-01-06	0.44	48	39	1	0.31

Split dataset into train and test

```
In [43]: # Split dataset into train and test
start1 = pd.datetime(1950, 1, 1)
end1 = pd.datetime(2009, 12, 31)
start2 = pd.datetime(2010,1,1)
end2 = pd.datetime(2017,12,31)

df_train = df[(df.DATE >= start1) & (df.DATE <= end1)]
df_test = df[(df.DATE >= start2) & (df.DATE <= end2)]
```

```
In [44]: # Drop all features excluding TMAX and average_PRCP. Can't use PRCP to predict future rain days.
feature_cols = ['TMAX','average_PRCP']

X = df_train[feature_cols]
y = df_train['RAIN']
X_test = df_test[feature_cols]
y_test = df_test['RAIN']
```

Null accuracy

```
In [45]: # calculate null accuracy
# Future models should be better than this, otherwise don't waste your time.
dumb = DummyClassifier(strategy='most_frequent')
dumb.fit(X, y)
y_dumb = dumb.predict(X_test)
print('Null accuracy:',metrics.accuracy_score(y_test, y_dumb))
```

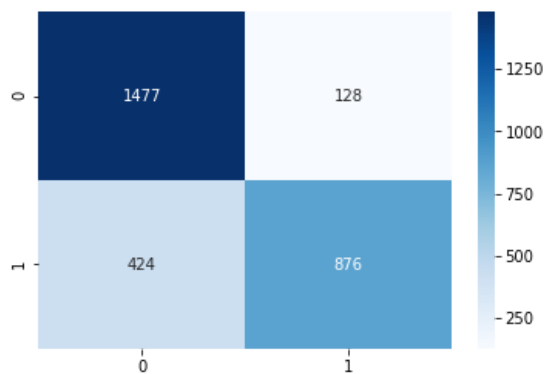
Null accuracy: 0.552495697074

Logistic regression model

```
In [48]: logmodel = LogisticRegression()
logmodel.fit(X,y)
y_pred = logmodel.predict(X_test)
print ('Accuracy:{:.2}'.format(metrics.accuracy_score(y_test, y_pred)))
sns.heatmap(confusion_matrix(y_test, y_pred), cmap='Blues', annot=True, fmt='g')
```

Accuracy:0.81

Out[48]: <matplotlib.axes._subplots.AxesSubplot at 0x10f686390>



```
In [49]: print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.78	0.92	0.84	1605
1	0.87	0.67	0.76	1300
avg / total	0.82	0.81	0.81	2905

```
In [50]: # Predicting using train set
y_pred = logmodel.predict(X)
print('Accuracy: {:.2}'.format(metrics.accuracy_score(y,y_pred)))
```

Accuracy: 0.81

```
In [51]: # Predicting using test set
y_pred = logmodel.predict(X_test)
print('Accuracy: {:.2}'.format(metrics.accuracy_score(y_test,y_pred)))
```

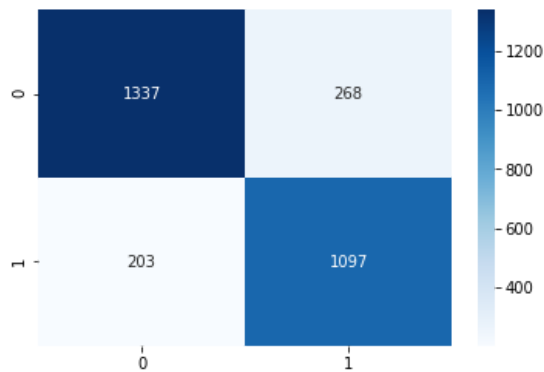
Accuracy: 0.81

K-Nearest Neighbour model

```
In [56]: knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X,y)
y_pred = knn.predict(X_test)
print('Accuracy:{:.2}'.format(metrics.accuracy_score(y_test, y_pred)))
sns.heatmap(confusion_matrix(y_test, y_pred), cmap='Blues', annot=True, fmt='g')
```

Accuracy:0.84

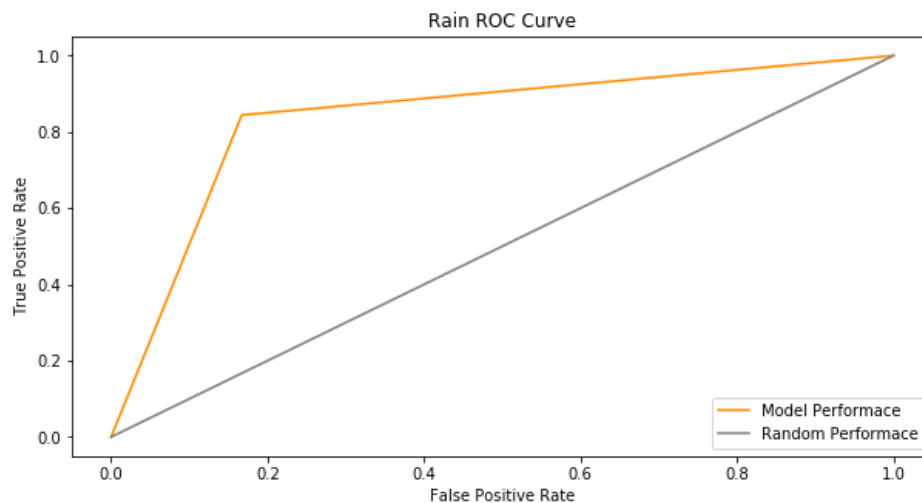
Out[56]: <matplotlib.axes._subplots.AxesSubplot at 0x1a1932bcc0>



```
In [68]: fpr, tpr, _ = metrics.roc_curve(y_test, y_pred)

fig, ax = plt.subplots(1, figsize=(10, 5))
plt.plot(fpr, tpr, color='darkorange', label='Model Performace')
plt.plot([0, 1], [0, 1], color='gray', label='Random Performace')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Rain ROC Curve')
plt.legend(loc="lower right")
```

Out[68]: <matplotlib.legend.Legend at 0x1a1ec7be80>



```
In [62]: print(classification_report(y_test,y_pred))
fpr, tpr, thresholds = roc_curve(y_test, y_pred)
```

	precision	recall	f1-score	support
0	0.87	0.83	0.85	1605
1	0.80	0.84	0.82	1300
avg / total	0.84	0.84	0.84	2905

```
In [63]: # Predicting using train set
y_pred = knn.predict(X)
print('Accuracy: {:.2}'.format(metrics.accuracy_score(y,y_pred)))
```

Accuracy: 0.85

```
In [64]: # Predicting using test set
y_pred = knn.predict(X_test)
print('Accuracy: {:.2}'.format(metrics.accuracy_score(y_test,y_pred)))
```

Accuracy: 0.84