



Final Project: Classification with Python

Table of Contents

- [Instructions](#)
- [About the Data](#)
- [Importing Data](#)
- [Data Preprocessing](#)
- [One Hot Encoding](#)
- [Train and Test Data Split](#)
- [Train Logistic Regression, KNN, Decision Tree, SVM, and Linear Regression models and return their appropriate accuracy scores](#)

Estimated Time Needed: **180 min**

Instructions

In this notebook, you will practice all the classification algorithms that we have learned in this course.

Below, is where we are going to use the classification algorithms to create a model based on our training data and evaluate our testing data using evaluation metrics learned in the course.

We will use some of the algorithms taught in the course, specifically:

1. Linear Regression
2. KNN
3. Decision Trees
4. Logistic Regression
5. SVM

We will evaluate our models using:

1. Accuracy Score
2. Jaccard Index
3. F1-Score
4. LogLoss

5. Mean Absolute Error
6. Mean Squared Error
7. R2-Score

Finally, you will use your models to generate the report at the end.

About The Dataset

The original source of the data is Australian Government's Bureau of Meteorology and the latest data can be gathered from <http://www.bom.gov.au/climate/dwo/>.

The dataset to be used has extra columns like 'RainToday' and our target is 'RainTomorrow', which was gathered from the Rattle at <https://bitbucket.org/kayontoga/rattle/src/master/data/weatherAUS.RData>

This dataset contains observations of weather metrics for each day from 2008 to 2017. The **weatherAUS.csv** dataset includes the following fields:

Field	Description	Unit	Type
Date	Date of the Observation in YYYY-MM-DD	Date	object
Location	Location of the Observation	Location	object
MinTemp	Minimum temperature	Celsius	float
MaxTemp	Maximum temperature	Celsius	float
Rainfall	Amount of rainfall	Millimeters	float
Evaporation	Amount of evaporation	Millimeters	float
Sunshine	Amount of bright sunshine	hours	float
WindGustDir	Direction of the strongest gust	Compass Points	object
WindGustSpeed	Speed of the strongest gust	Kilometers/ Hour	object
WindDir9am	Wind direction averaged of 10 minutes prior to 9am	Compass Points	object
WindDir3pm	Wind direction averaged of 10 minutes prior to 3pm	Compass Points	object
WindSpeed9am	Wind speed averaged of 10 minutes prior to 9am	Kilometers/ Hour	float
WindSpeed3pm	Wind speed averaged of 10 minutes prior to 3pm	Kilometers/ Hour	float
Humidity9am	Humidity at 9am	Percent	float
Humidity3pm	Humidity at 3pm	Percent	float
Pressure9am	Atmospheric pressure reduced to mean sea level at 9am	Hectopascal	float

Field	Description	Unit	Type
Pressure3pm	Atmospheric pressure reduced to mean sea level at 3pm	Hectopascal	float
Cloud9am	Fraction of the sky obscured by cloud at 9am	Eights	float
Cloud3pm	Fraction of the sky obscured by cloud at 3pm	Eights	float
Temp9am	Temperature at 9am	Celsius	float
Temp3pm	Temperature at 3pm	Celsius	float
RainToday	If there was rain today	Yes/No	object
RainTomorrow	If there is rain tomorrow	Yes/No	float

Column definitions were gathered from <http://www.bom.gov.au/climate/dwo/IDCJDW0000.shtml>

Import the required libraries

```
In [ ]: # All Libraries required for this lab are listed below. The libraries pre
!pip install pandas numpy seaborn matplotlib scikit-learn
# Note: If your environment doesn't support "!mamba install", use "!pip i
```

```
In [ ]: # Surpress warnings:
def warn(*args, **kwargs):
    pass
import warnings
warnings.warn = warn
```

```
In [2]: import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import LinearRegression
from sklearn import preprocessing
import numpy as np
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn import svm
from sklearn.metrics import jaccard_score
from sklearn.metrics import f1_score
from sklearn.metrics import log_loss
from sklearn.metrics import confusion_matrix, accuracy_score
import sklearn.metrics as metrics
```

Importing the Dataset

```
In [ ]: import requests

def download(url, filename):
    response = requests.get(url)
    if response.status_code == 200:
        with open(filename, "wb") as f:
```

```
f.write(response.content)
```

```
In [ ]: path='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/'
```

```
In [ ]: download(path, "Weather_Data.csv")
filename = "Weather_Data.csv"
```

```
In [3]: df = pd.read_csv("Weather_Data.csv")
```

Note: This version of the lab is designed for JupyterLite, which necessitates downloading the dataset to the interface. However, when working with the downloaded version of this notebook on your local machines (Jupyter Anaconda), you can simply **skip the steps above of "Importing the Dataset"** and use the URL directly in the `pandas.read_csv()` function. You can uncomment and run the statements in the cell below.

```
In [ ]: #filepath = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain
#df = pd.read_csv(filepath)
```

```
In [4]: df.head()
```

```
Out[4]:
```

	Date	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGu
--	------	---------	---------	----------	-------------	----------	-------------	--------

0	2/1/2008	19.5	22.4	15.6	6.2	0.0	W
1	2/2/2008	19.5	25.6	6.0	3.4	2.7	W
2	2/3/2008	21.6	24.5	6.6	2.4	0.1	W
3	2/4/2008	20.2	22.8	18.8	2.2	0.0	W
4	2/5/2008	19.7	25.7	77.4	4.8	0.0	W

5 rows × 22 columns

Data Preprocessing

One Hot Encoding

First, we need to perform one hot encoding to convert categorical variables to binary variables.

```
In [5]: df_sydney_processed = pd.get_dummies(data=df, columns=['RainToday', 'Wind
```

Next, we replace the values of the 'RainTomorrow' column changing them from a categorical column to a binary column. We do not use the `get_dummies` method because we would end up with two columns for 'RainTomorrow' and we do not want, since 'RainTomorrow' is our target.

```
In [6]: df_sydney_processed.replace(['No', 'Yes'], [0,1], inplace=True)
```

Training Data and Test Data

Now, we set our 'features' or x values and our Y or target variable.

```
In [7]: df_sydney_processed.drop('Date',axis=1,inplace=True)
```

```
In [8]: df_sydney_processed = df_sydney_processed.astype(float)
```

```
In [9]: features = df_sydney_processed.drop(columns='RainTomorrow', axis=1)
Y = df_sydney_processed['RainTomorrow']
```

Linear Regression

Q1) Use the `train_test_split` function to split the features and Y dataframes with a `test_size` of 0.2 and the `random_state` set to 10.

```
In [ ]: #Enter Your Code and Execute
```

```
In [10]: x_train, x_test, y_train, y_test = train_test_split(features, Y, test_size=0.2, random_state=10)
```

Q2) Create and train a Linear Regression model called `LinearReg` using the training data (`x_train`, `y_train`).

```
In [ ]: #Enter Your Code and Execute
```

```
In [11]: LinearReg = LinearRegression()
LinearReg.fit(x_train, y_train)
```

```
Out[11]: ▼ LinearRegression ⓘ ?
LinearRegression()
```

Q3) Now use the `predict` method on the testing data (`x_test`) and save it to the array `predictions`.

```
In [ ]: #Enter Your Code and Execute
```

```
In [12]: predictions = LinearReg.predict(x_test)
```

Q4) Using the `predictions` and the `y_test` dataframe calculate the value for each metric using the appropriate function.

```
In [ ]: #Enter Your Code and Execute
```

```
In [13]: LinearRegression_MAE = metrics.mean_absolute_error(y_test, predictions)
LinearRegression_MSE = metrics.mean_squared_error(y_test, predictions)
LinearRegression_R2 = metrics.r2_score(y_test, predictions)
```

Q5) Show the MAE, MSE, and R2 in a tabular format using data frame for the linear model.

In []: *#Enter Your Code and Execute*

```
In [14]: Report = pd.DataFrame({"Model": ["Linear Regression"],
                                "MAE": [LinearRegression_MAE],
                                "MSE": [LinearRegression_MSE],
                                "R^2": [LinearRegression_R2]}).set_index
                                Report
```

```
Out[14]:
```

	MAE	MSE	R^2
Model			
Linear Regression	0.265804	0.129578	0.309619

KNN

Q6) Create and train a KNN model called KNN using the training data (x_train, y_train) with the n_neighbors parameter set to 4.

In []: *#Enter Your Code and Execute*

```
In [15]: KNN = KNeighborsClassifier(n_neighbors=4)
KNN.fit(x_train, y_train)
```

```
Out[15]:
```

▼ KNeighborsClassifier ⓘ ?

KNeighborsClassifier(n_neighbors=4)

Q7) Now use the predict method on the testing data (x_test) and save it to the array predictions.

In []: *#Enter Your Code and Execute*

```
In [16]: predictions = KNN.predict(x_test)
```

Q8) Using the predictions and the y_test dataframe calculate the value for each metric using the appropriate function.

In []: *#Enter Your Code and Execute*

```
In [17]: KNN_Accuracy_Score = metrics.accuracy_score(y_test, predictions)
KNN_JaccardIndex = metrics.jaccard_score(y_test, predictions)
KNN_F1_Score = metrics.f1_score(y_test, predictions)
```

Decision Tree

Q9) Create and train a Decision Tree model called Tree using the training

```
data(x_train, y_train).
```

```
In [ ]: #Enter Your Code and Execute
```

```
In [18]: Tree = DecisionTreeClassifier()  
Tree.fit(x_train, y_train)
```

```
Out[18]: ▼ DecisionTreeClassifier ⓘ ?  
DecisionTreeClassifier()
```

Q10) Now use the `predict` method on the testing data (`x_test`) and save it to the array `predictions`.

```
In [ ]: #Enter Your Code and Execute
```

```
In [19]: predictions = Tree.predict(x_test)
```

Q11) Using the `predictions` and the `y_test` dataframe calculate the value for each metric using the appropriate function.

```
In [ ]: #Enter Your Code and Execute
```

```
In [20]: Tree_Accuracy_Score = metrics.accuracy_score(y_test, predictions)  
Tree_JaccardIndex = metrics.jaccard_score(y_test, predictions)  
Tree_F1_Score = metrics.f1_score(y_test, predictions)
```

Logistic Regression

Q12) Use the `train_test_split` function to split the features and `Y` dataframes with a `test_size` of `0.2` and the `random_state` set to `1`.

```
In [ ]: #Enter Your Code and Execute
```

```
In [21]: x_train, x_test, y_train, y_test = train_test_split(features, Y, test_siz
```

Q13) Create and train a LogisticRegression model called `LR` using the training data (`x_train`, `y_train`) with the `solver` parameter set to `liblinear`.

```
In [ ]: #Enter Your Code and Execute
```

```
In [22]: LR = LogisticRegression(solver='liblinear')  
LR.fit(x_train, y_train)
```

```
Out[22]: ▼ LogisticRegression ⓘ ?  
LogisticRegression(solver='liblinear')
```

Q14) Now, use the `predict` and `predict_proba` methods on the testing data (`x_test`) and save it as 2 arrays `predictions` and `predict_proba`.

```
In [ ]: #Enter Your Code and Execute
```

```
In [23]: predictions = LR.predict(x_test)
```

```
In [24]: predict_proba = LR.predict_proba(x_test)
```

Q15) Using the `predictions`, `predict_proba` and the `y_test` dataframe calculate the value for each metric using the appropriate function.

```
In [ ]: #Enter Your Code and Execute
```

```
In [25]: LR_Accuracy_Score = metrics.accuracy_score(y_test, predictions)
LR_JaccardIndex = metrics.jaccard_score(y_test, predictions)
LR_F1_Score = metrics.f1_score(y_test, predictions)
LR_Log_Loss = metrics.log_loss(y_test, predict_proba)
```

SVM

Q16) Create and train a SVM model called SVM using the training data (`x_train`, `y_train`).

```
In [ ]: #Enter Your Code and Execute
```

```
In [26]: SVM = SVC()
SVM.fit(x_train, y_train)
```

```
Out[26]: ▼ SVC ⓘ ?
SVC()
```

Q17) Now use the `predict` method on the testing data (`x_test`) and save it to the array `predictions`.

```
In [ ]: #Enter Your Code and Execute
```

```
In [27]: predictions = SVM.predict(x_test)
```

Q18) Using the `predictions` and the `y_test` dataframe calculate the value for each metric using the appropriate function.

```
In [28]: SVM_Accuracy_Score = metrics.accuracy_score(y_test, predictions)
SVM_JaccardIndex = metrics.jaccard_score(y_test, predictions)
SVM_F1_Score = metrics.f1_score(y_test, predictions)
```


Report

Q19) Show the Accuracy, Jaccard Index, F1-Score and LogLoss in a tabular format using data frame for all of the above models.

*LogLoss is only for Logistic Regression Model

```
In [29]: Report = pd.DataFrame({"Model": ["KNN", "Decision Tree", "Logistic Regres",
                                         "Accuracy": [KNN_Accuracy_Score, Tree_Ac
                                         "Jaccard Index": [KNN_JaccardIndex, Tree
                                         "F1 Score": [KNN_F1_Score, Tree_F1_Score
                                         "Log Loss": [np.nan, np.nan, LR_Log_Loss

Report
```

```
Out[29]:
```

	Accuracy	Jaccard Index	F1 Score	Log Loss
Model				
KNN	0.795420	0.309278	0.472441	NaN
Decision Tree	0.743511	0.325301	0.490909	NaN
Logistic Regression	0.836641	0.509174	0.674772	0.381259
SVM	0.722137	0.000000	0.000000	NaN

How to submit

Once you complete your notebook you will have to share it. You can download the notebook by navigating to "File" and clicking on "Download" button.

This will save the (.ipynb) file on your computer. Once saved, you can upload this file in the "My Submission" tab, of the "Peer-graded Assignment" section.

About the Authors:

[Joseph Santarcangelo](#) has a PhD in Electrical Engineering, his research focused on using machine learning, signal processing, and computer vision to determine how videos impact human cognition. Joseph has been working for IBM since he completed his PhD.

Other Contributors

[Svitlana Kramar](#)

© IBM Corporation 2020. All rights reserved.

```
<!-- ## Change Log | Date (YYYY-MM-DD) | Version | Changed By | Change Description | |
----- | ----- | ----- | ----- | 2022-06-22 | 2.0 | Svitlana K. |
Deleted GridSearch and Mock | --!>
```

