



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

Field	Description	Unit	Type
Pressure9am	Atmospheric pressure reduced to mean sea level at 9am	Hectopascal	float
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-installed
# !mamba install -qy pandas==1.3.4 numpy==1.21.4 seaborn==0.9.0 matplotlib==3.5.0 s
# Note: If your environment doesn't support "!mamba install", use "!pip install"
```

```
In [125...]: # Suppress warnings:
def warn(*args, **kwargs):
    pass
import warnings
warnings.warn = warn
```

```
In [126...]: #you are running the Lab in your browser, so we will install the Libraries using -
import piplite
await piplite.install(['pandas'])
await piplite.install(['numpy'])
```

```
In [127...]: 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 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 [128...]:  
from pyodide.http import pyfetch  
  
async def download(url, filename):  
    response = await pyfetch(url)  
    if response.status == 200:  
        with open(filename, "wb") as f:  
            f.write(await response.bytes())  
  
In [129...]: path='https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-ML0101EN-SkillsNetwork/labs/Module%203/Weather_Datasource.csv'  
  
In [130...]: await download(path, "Weather_Data.csv")  
filename ="Weather_Data.csv"  
  
In [131...]: df = pd.read_csv("Weather_Data.csv")  
df.head()  
  
Out[131]:  


|   | Date     | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine | WindGustDir | WindGus |
|---|----------|---------|---------|----------|-------------|----------|-------------|---------|
| 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 [132...]: df_sydney_processed = pd.get_dummies(data=df, columns=['RainToday', 'WindGustDir',
```

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 [133]: 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 [134]: df_sydney_processed.drop('Date', axis=1, inplace=True)
```

```
In [49]: df_sydney_processed = df_sydney_processed.astype(float)
```

```
In [135]: 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 [51]: #Enter Your Code, Execute and take the Screenshot
```

```
In [136]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split( features, Y, test_size=0.2, random_state=10)
print ('Train set:', X_train.shape, y_train.shape)
print ('Test set:', X_test.shape, y_test.shape)
```

Train set: (2616, 66) (2616,)

Test set: (655, 66) (655,)

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

```
In [18]: #Enter Your Code, Execute and take the Screenshot
```

```
In [137]: from sklearn import linear_model
LinearReg = linear_model.LinearRegression()
LinearReg.fit(X_train,y_train)
print('Coefficients:', LinearReg.coef_)
print('Intercept:', LinearReg.intercept_)
```

```
Coefficients: [-2.36883601e-02  1.30031554e-02  7.30154915e-04  6.49290860e-03
 -3.51614515e-02  4.23841910e-03  1.82843563e-03  7.89794445e-04
  9.55888360e-04  8.56007751e-03  7.70261345e-03  -9.24932344e-03
 -8.87528023e-03  1.00457553e-02  1.44651992e-02  -3.48325241e-03
  1.03190076e+10  1.03190076e+10  -4.35428535e+09  -4.35428535e+09
 -4.35428535e+09  -4.35428535e+09  -4.35428535e+09  -4.35428535e+09
 -4.35428535e+09  -4.35428535e+09  -4.35428535e+09  -4.35428535e+09
 -4.35428535e+09  -4.35428535e+09  -4.35428535e+09  -4.35428535e+09
 -4.35428535e+09  -4.35428535e+09  4.72755310e+09  4.72755310e+09
  4.72755310e+09  4.72755310e+09  4.72755310e+09  4.72755310e+09
  4.72755310e+09  4.72755310e+09  4.72755310e+09  4.72755310e+09
  4.72755310e+09  4.72755310e+09  -1.18315500e+10  -1.18315500e+10
 -1.18315500e+10  -1.18315500e+10  -1.18315500e+10  -1.18315500e+10
 -1.18315500e+10  -1.18315500e+10  -1.18315500e+10  -1.18315500e+10
 -1.18315500e+10  -1.18315500e+10]

Intercept: 1139274582.2150812
```

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

In []: #Enter Your Code, Execute and take the Screenshot

In [138...]: 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, Execute and take the Screenshot

In [139...]:
`from sklearn.metrics import r2_score`
`LinearRegression_MAE = np.mean(np.absolute(y_test-y_test))`
`LinearRegression_MSE = np.mean(np.absolute(y_test-y_test)**2)`
`LinearRegression_R2 = r2_score(y_test, y_test)`

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

In []: #Enter Your Code, Execute and take the Screenshot

In [140...]:
`Report = pd.DataFrame({`
 `'Metric': ['Mean Absolute Error (MAE)', 'Mean Squared Error (MSE)', 'R-squared'],`
 `'Value': [LinearRegression_MAE, LinearRegression_MSE, LinearRegression_R2]})`
`print(Report)`

	Metric	Value
0	Mean Absolute Error (MAE)	0.0
1	Mean Squared Error (MSE)	0.0
2	R-squared (R2)	1.0

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 Below, Execute, and Save the Screenshot of the Final Output

In [141...]: `KNN = KNeighborsClassifier(n_neighbors=4).fit(X_train,y_train)`
KNN

Out[141]:

```
▼      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 Below, Execute, and Save the Screenshot of the Final Output

In [142...]: `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 Below, Execute, and Save the Screenshot of the Final Output

In [143...]: `KNN_Accuracy_Score = metrics.accuracy_score(y_test, predictions)`
`KNN_JaccardIndex = jaccard_score(y_test, predictions, pos_label=0)`
`KNN_F1_Score = f1_score(y_test, predictions)`
`print('KNN Accuracy:', KNN_Accuracy_Score)`
`print('KNN Jaccard:', KNN_JaccardIndex)`
`print('KNN F1:', KNN_F1_Score)`

KNN Accuracy: 0.8183206106870229

KNN Jaccard: 0.7901234567901234

KNN F1: 0.5966101694915255

Decision Tree

Q9) Create and train a Decision Tree model called Tree using the training data (`x_train` , `y_train`).

In []: #Enter Your Code, Execute and take the Screenshot

In [144...]: `Tree = DecisionTreeClassifier(criterion="entropy")`
`Tree.fit(X_train,y_train)`

Out[144]:

```
▼ DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy')
```

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

In []: #Enter Your Code, Execute and take the Screenshot

In [145...]: 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, Execute and take the Screenshot

```
In [146...]: Tree_Accuracy_Score = metrics.accuracy_score(y_test, predictions)
Tree_JaccardIndex = jaccard_score(y_test, predictions, pos_label=0)
Tree_F1_Score = f1_score(y_test, predictions)
print('Tree Accuracy:', Tree_Accuracy_Score)
print('Tree Jaccard:', Tree_JaccardIndex)
print('Tree F1:', Tree_F1_Score)
```

Tree Accuracy: 0.76793893129771
 Tree Jaccard: 0.7185185185186
 Tree F1: 0.6020942408376962

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, Execute and take the Screenshot

```
In [147...]: x_train, x_test, y_train, y_test = train_test_split( features, Y, test_size=0.2, random_state=1)
print ('Train set:', X_train.shape, y_train.shape)
print ('Test set:', X_test.shape, y_test.shape)
```

Train set: (2616, 66) (2616,)
 Test set: (655, 66) (655,)

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, Execute and take the Screenshot

In [148...]: LR = LogisticRegression(C=0.01, solver='liblinear').fit(X_train,y_train)

```
LR
Out[148]: LogisticRegression
LogisticRegression(C=0.01, 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, Execute and take the Screenshot
In [149... predictions = LR.predict(X_test)
In [150... predict_proba = LR.predict_proba(X_test)
predict_proba
Out[150]: array([[0.64085492, 0.35914508],
       [0.75061354, 0.24938646],
       [0.70553245, 0.29446755],
       ...,
       [0.7453308 , 0.2546692 ],
       [0.74595521, 0.25404479],
       [0.79379072, 0.20620928]])
```

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, Execute and take the Screenshot
In [151... LR_Accuracy_Score = metrics.accuracy_score(y_test, predictions)
LR_JaccardIndex = jaccard_score(y_test, predictions, pos_label=0)
LR_F1_Score = f1_score(y_test, predictions)
LR_Log_Loss = log_loss(y_test, predict_proba)
print('LR Accuracy:', LR_Accuracy_Score)
print('LR Jaccard:', LR_JaccardIndex)
print('LR F1:', LR_F1_Score)
print('LR Log Loss:', LR_Log_Loss)
LR Accuracy: 0.7221374045801526
LR Jaccard: 0.7221374045801526
LR F1: 0.0
LR Log Loss: 0.5971474619677751
```

SVM

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

In []: #Enter Your Code Below, Execute, and Save the Screenshot of the Final Output

In [152...]

```
from sklearn import svm
SVM= svm.SVC(kernel='rbf')
SVM.fit(X_train, y_train)
```

Out[152]:

```
▼ 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 Below, Execute, and Save the Screenshot of the Final Output

In [155...]

```
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 [156...]

```
SVM_Accuracy_Score = metrics.accuracy_score(y_test, predictions)
SVM_JaccardIndex = jaccard_score(y_test, predictions, pos_label=0)
SVM_F1_Score = f1_score(y_test, predictions)
print('SVM Accuracy:', SVM_Accuracy_Score)
print('SVM Jaccard:', SVM_JaccardIndex)
print('SVM F1:', SVM_F1_Score)
```

SVM Accuracy: 0.7221374045801526

SVM Jaccard: 0.7221374045801526

SVM F1: 0.0

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 [158...]

```
models_metrics = {
    'Model': ['KNN', 'Decision Tree', 'SVM', 'Logistic Regression'],
    'Accuracy': [KNN_Accuracy_Score, Tree_Accuracy_Score, SVM_Accuracy_Score, LR_Accuracy],
    'Jaccard Index': [KNN_JaccardIndex, Tree_JaccardIndex, SVM_JaccardIndex, LR_Jaccard_Index],
    'F1-Score': [KNN_F1_Score, Tree_F1_Score, SVM_F1_Score, LR_F1_Score],
    'LogLoss': [np.nan, np.nan, np.nan, LR_Log_Loss]
}
Report = pd.DataFrame(models_metrics)

# Display the DataFrame
print(Report)
```

	Model	Accuracy	Jaccard Index	F1-Score	LogLoss
0	KNN	0.818321	0.790123	0.596610	NaN
1	Decision Tree	0.767939	0.718519	0.602094	NaN
2	SVM	0.722137	0.722137	0.000000	NaN
3	Logistic Regression	0.722137	0.722137	0.000000	0.597147

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](#)

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2022-06-22	2.0	Svitlana K.	Deleted GridSearch and Mock

© IBM Corporation 2020. All rights reserved.