# Skills
# Network

# Final Project: Classification with Python

## Table of Contents

Estimated Time Needed: **180 min**

# Instructions

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

After completing this notebook, you will need to upload it to the "Submit Your Work and Review Your Peers" section of the Final Project module.

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 displaying the accuracy scores.

# 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 |
| 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 |
| RISK_MM | Amount of rain tomorrow | Millimeters | float |
| RainTomorrow | If there is rain tomorrow | Yes/No | float |

Column definitions were gathered from

## Import the required libraries

```
In [1]:    # All Libraries required for this lab are listed below. The libraries pre-installed on S
           # !mamba install -qy pandas==1.3.4 numpy==1.21.4 seaborn==0.9.0 matplotlib==3.5.0 scikit
           # Note: If your environment doesn't support "!mamba install", use "!pip install"
```

```
In [2]:    # Surpress warnings:
           def warn(*args, **kwargs):
               pass
           import warnings
           warnings.warn = warn
```

```
In [3]:    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 GridSearchCV
           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
           import matplotlib.pyplot as plt
           from sklearn.metrics import confusion_matrix, accuracy_score
           import sklearn.metrics as metrics
```

### Importing the Dataset

```
In [4]:    df = pd.read_csv('Weather_Data.csv')

           df.head()
```

Out[4]:

| | Date | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine | WindGustDir | WindGustSpeed | WindDir9am | |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2/1/2008 | 19.5 | 22.4 | 15.6 | 6.2 | 0.0 | W | 41 | S | |
| 1 | 2/2/2008 | 19.5 | 25.6 | 6.0 | 3.4 | 2.7 | W | 41 | W | |
| 2 | 2/3/2008 | 21.6 | 24.5 | 6.6 | 2.4 | 0.1 | W | 41 | ESE | |
| 3 | 2/4/2008 | 20.2 | 22.8 | 18.8 | 2.2 | 0.0 | W | 41 | NNE | |
| 4 | 2/5/2008 | 19.7 | 25.7 | 77.4 | 4.8 | 0.0 | W | 41 | NNE | |

5 rows × 22 columns

## Data Preprocessing

### Transforming Categorical Variables

First, we need to convert categorical variables to binary variables. We will use pandas `get_dummies()` method for this.

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

```
In [7]:  df_sydney_processed.head()
```

Out[7]:

| | Date | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine | WindGustSpeed | WindSpeed9am | WindSpeed3 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2/1/2008 | 19.5 | 22.4 | 15.6 | 6.2 | 0.0 | 41 | 17 | |
| 1 | 2/2/2008 | 19.5 | 25.6 | 6.0 | 3.4 | 2.7 | 41 | 9 | |
| 2 | 2/3/2008 | 21.6 | 24.5 | 6.6 | 2.4 | 0.1 | 41 | 17 | |
| 3 | 2/4/2008 | 20.2 | 22.8 | 18.8 | 2.2 | 0.0 | 41 | 22 | |
| 4 | 2/5/2008 | 19.7 | 25.7 | 77.4 | 4.8 | 0.0 | 41 | 11 | |

5 rows × 68 columns

## Training Data and Test Data

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

```
In [8]:  df_sydney_processed.drop('Date',axis=1,inplace=True)
```

```
In [9]:  df_sydney_processed = df_sydney_processed.astype(float)
```

```
In [10]:  features = df_sydney_processed.drop(columns='RainTomorrow', axis=1)
          Y = df_sydney_processed['RainTomorrow']
```

```
In [11]:  features
```

Out[11]:

| | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine | WindGustSpeed | WindSpeed9am | WindSpeed3pm | H |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 19.5 | 22.4 | 15.6 | 6.2 | 0.0 | 41.0 | 17.0 | 20.0 | |
| 1 | 19.5 | 25.6 | 6.0 | 3.4 | 2.7 | 41.0 | 9.0 | 13.0 | |
| 2 | 21.6 | 24.5 | 6.6 | 2.4 | 0.1 | 41.0 | 17.0 | 2.0 | |
| 3 | 20.2 | 22.8 | 18.8 | 2.2 | 0.0 | 41.0 | 22.0 | 20.0 | |
| 4 | 19.7 | 25.7 | 77.4 | 4.8 | 0.0 | 41.0 | 11.0 | 6.0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 3266 | 8.6 | 19.6 | 0.0 | 2.0 | 7.8 | 37.0 | 22.0 | 20.0 | |
| 3267 | 9.3 | 19.2 | 0.0 | 2.0 | 9.2 | 30.0 | 20.0 | 7.0 | |
| 3268 | 9.4 | 17.7 | 0.0 | 2.4 | 2.7 | 24.0 | 15.0 | 13.0 | |
| 3269 | 10.1 | 19.3 | 0.0 | 1.4 | 9.3 | 43.0 | 17.0 | 19.0 | |
| 3270 | 7.6 | 19.3 | 0.0 | 3.4 | 9.4 | 35.0 | 13.0 | 13.0 | |

3271 rows × 66 columns

```
In [12]:  Y

Out[12]:  0         1.0
          1         1.0
          2         1.0
          3         1.0
          4         1.0
                   ...
          3266      0.0
          3267      0.0
          3268      0.0
          3269      0.0
          3270      0.0
          Name: RainTomorrow, Length: 3271, dtype: float64
```

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

```
In [14]:  x_train, x_test, y_train, y_test = train_test_split(features,Y,test_size=0.2,  random_sta
```

```
In [15]:  x_train.values, x_test.values
```

```
Out[15]:  (array([[14.8, 22. , 33.8, ...,  0. ,  0. ,  0. ],
                  [ 8.1, 18.4,  0. , ...,  0. ,  0. ,  0. ],
                  [15.4, 21.1,  0. , ...,  1. ,  0. ,  0. ],
                  ...,
                  [ 6.7, 17.3,  0. , ...,  0. ,  0. ,  0. ],
                  [15. , 22.7,  9.4, ...,  0. ,  0. ,  0. ],
                  [15.9, 30.1,  0.2, ...,  0. ,  0. ,  0. ]]),
           array([[18.7, 35.7,  0.2, ...,  0. ,  0. ,  0. ],
                  [15. , 24.8, 22.4, ...,  0. ,  0. ,  0. ],
                  [ 8.5, 16. ,  3.4, ...,  0. ,  0. ,  0. ],
                  ...,
                  [11.6, 15.7, 18.2, ...,  0. ,  0. ,  0. ],
                  [16.9, 23.9,  0. , ...,  0. ,  0. ,  0. ],
                  [15.6, 22.2, 22.8, ...,  0. ,  0. ,  0. ]]))
```

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

```
In [16]:  #Enter Your Code, Execute and take the Screenshot
```

```
In [17]:  LinearReg = LinearRegression()
```

```
In [18]:  LinearReg.fit(x_train, y_train)
```

```
Out[18]:  LinearRegression()
```

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

```
In [19]:  #Enter Your Code, Execute and take the Screenshot
```

```
In [20]:   predictions = LinearReg.predict(x_test)
```

```
In [21]:   predictions[:5]
```

```
Out[21]:   array([0.13184071, 0.2761859 , 0.97818819, 0.2874561 , 0.13241371])
```

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

```
In [22]:   #Enter Your Code, Execute and take the Screenshot

           from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
```

```
In [23]:   LinearRegression_MAE = mean_absolute_error(y_test,predictions)
           LinearRegression_MSE = mean_squared_error(y_test,predictions)
           LinearRegression_R2 = r2_score(y_test,predictions)
```

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

```
In [24]:   #Enter Your Code, Execute and take the Screenshot

           d = {'MAE': LinearRegression_MAE, 'MSE': LinearRegression_MSE, 'R2': LinearRegression_R2
```

```
In [25]:   LRReport = pd.DataFrame(d, index=[0])
```

```
In [26]:   LRReport
```

Out[26]:

|   | MAE | MSE | R2 | Loss |
|---|---|---|---|---|
| **0** | 0.256318 | 0.115721 | 0.427132 | NaN |

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

```
In [28]:   KNN = KNeighborsClassifier(n_neighbors=4)
```

```
In [29]:   KNN.fit(x_train,y_train)
```

```
Out[29]:   KNeighborsClassifier(n_neighbors=4)
```

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

```
In [30]:   #Enter Your Code, Execute and take the Screenshot
```

```
In [31]:   predictions = KNN.predict(x_test)
```

```
In [32]:   predictions[:5]
```

```
Out[32]:    array([0., 0., 1., 0., 0.])
```

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

```
In [33]:   #Enter Your Code, Execute and take the Screenshot
```

```
In [34]:   KNN_Accuracy_Score = accuracy_score(y_test,predictions)
           KNN_JaccardIndex = jaccard_score(y_test,predictions)
           KNN_F1_Score = f1_score(y_test,predictions)
```

```
In [35]:   #Enter Your Code, Execute and take the Screenshot

           d = {'Accuracy': KNN_Accuracy_Score, 'Jaccard': KNN_JaccardIndex, 'F1': KNN_F1_Score, 'L
```

```
In [36]:   KNNReport = pd.DataFrame(d, index=[1])
```

```
In [37]:   KNNReport
```

```
Out[37]:
```

|   | Accuracy | Jaccard | F1 | Loss |
|---|---|---|---|---|
| 1 | 0.818321 | 0.425121 | 0.59661 | NaN |

## Decision Tree

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

```
In [38]:   #Enter Your Code, Execute and take the Screenshot
```

```
In [39]:   Tree = DecisionTreeClassifier()
```

```
In [40]:   Tree.fit(x_train,y_train)
```

```
Out[40]:   DecisionTreeClassifier()
```

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

```
In [41]:   #Enter Your Code, Execute and take the Screenshot
```

```
In [42]:   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 [43]:   Tree_Accuracy_Score = accuracy_score(y_test,predictions)
           Tree_JaccardIndex = jaccard_score(y_test,predictions)
           Tree_F1_Score = f1_score(y_test,predictions)
```

```
In [44]:   #Enter Your Code, Execute and take the Screenshot

           d = {'Accuracy': Tree_Accuracy_Score, 'Jaccard': Tree_JaccardIndex, 'F1': Tree_F1_Score,
```

```
In [45]:  TreeReport = pd.DataFrame(d, index=[2])
```

```
In [46]:  TreeReport
```

Out[46]:

| | **Accuracy** | **Jaccard** | **F1** | **Loss** |
|---|---|---|---|---|
| **2** | 0.757252 | 0.395437 | 0.566757 | NaN |

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

```
In [48]:  x_train, x_test, y_train, y_test = train_test_split(features,Y,test_size=0.2, random_sta
```

**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 [49]:  #Enter Your Code, Execute and take the Screenshot
```

```
In [50]:  LR = LogisticRegression(solver='liblinear')
```

```
In [51]:  LR.fit(x_train,y_train)
```

Out[51]:  LogisticRegression(solver='liblinear')

**Q14) Now, use the `predict` method on the testing data ( `x_test` ) and save it to the array `predictions` .**

```
In [52]:  #Enter Your Code, Execute and take the Screenshot
```

```
In [53]:  predictions = LR.predict(x_test)
```

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

```
In [54]:  #Enter Your Code, Execute and take the Screenshot
```

```
In [55]:  LR_Accuracy_Score = accuracy_score(y_test,predictions)
          LR_JaccardIndex = jaccard_score(y_test,predictions)
          LR_F1_Score = f1_score(y_test,predictions)
          LR_Log_Loss = log_loss(y_test,predictions)
```

```
In [56]:  #Enter Your Code, Execute and take the Screenshot

          d = {'Accuracy': LR_Accuracy_Score, 'Jaccard': LR_JaccardIndex, 'F1': LR_F1_Score, 'Loss
```

```
In [57]:  LogReport = pd.DataFrame(d, index=[3])
```

```
In [58]:  LogReport
```

Out[58]:

| | Accuracy | Jaccard | F1 | Loss |
|---|---|---|---|---|
| 3 | 0.836641 | 0.509174 | 0.674772 | 5.642256 |

## SVM

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

In [59]: 
```
#Enter Your Code, Execute and take the Screenshot

from sklearn.svm import SVC
```

In [60]: 
```
SVM = SVC()
```

In [61]: 
```
SVM.fit(x_train,y_train)
```
Out[61]: 
```
SVC()
```

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

In [62]: 
```
#Enter Your Code, Execute and take the Screenshot
```

In [63]: 
```
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 [64]: 
```
SVM_Accuracy_Score = accuracy_score(y_test,predictions)
SVM_JaccardIndex = jaccard_score(y_test,predictions)
SVM_F1_Score = f1_score(y_test,predictions)
```

In [65]: 
```
#Enter Your Code, Execute and take the Screenshot

d = {'Accuracy': SVM_Accuracy_Score, 'Jaccard': SVM_JaccardIndex, 'F1': SVM_F1_Score, 'L
```

In [66]: 
```
SVMReport = pd.DataFrame(d, index=[4])
```

In [67]: 
```
SVMReport
```

Out[67]:

| | Accuracy | Jaccard | F1 | Loss |
|---|---|---|---|---|
| 4 | 0.722137 | 0.0 | 0.0 | NaN |

## 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 [68]: 
```
#Enter Your Code, Execute and take the Screenshot
```

```
In [69]:  Report = pd.concat([KNNReport, TreeReport, LogReport, SVMReport])

In [70]:  Report.rename(index={1: 'KNN', 2: 'Tree', 3: 'Logistic', 4: 'SVM'}, inplace=True)

In [71]:  Report
```

Out[71]:

|          | Accuracy | Jaccard  | F1       | Loss     |
|----------|----------|----------|----------|----------|
| KNN      | 0.818321 | 0.425121 | 0.596610 | NaN      |
| Tree     | 0.757252 | 0.395437 | 0.566757 | NaN      |
| Logistic | 0.836641 | 0.509174 | 0.674772 | 5.642256 |
| 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

Himanshu Birla

# Change Log

| Date (YYYY-MM-DD) | Version | Changed By    | Change Description          |
|-------------------|---------|---------------|-----------------------------|
| 2020-08-27        | 1.0     | Malika Singla | Added lab to GitLab         |
| 2022-06-22        | 2.0     | Lana K.       | Deleted GridSearch and Mock |