# Space X Falcon 9 First Stage Landing Prediction

## Assignment: Machine Learning Prediction

Estimated time needed: **60** minutes

Space X advertises Falcon 9 rocket launches on its website with a cost of 62 million dollars; other providers cost upward of 165 million dollars each, much of the savings is because Space X can reuse the first stage. Therefore if we can determine if the first stage will land, we can determine the cost of a launch. This information can be used if an alternate company wants to bid against space X for a rocket launch. In this lab, you will create a machine learning pipeline to predict if the first stage will land given the data from the preceding labs.



Several examples of an unsuccessful landing are shown here:

SEPTEMBER 2013   HARD IMPACT ON OCEAN

Most unsuccessful landings are planed. Space X; performs a controlled landing in the oceans.

# Objectives

Perform exploratory Data Analysis and determine Training Labels

- create a column for the class
- Standardize the data
- Split into training data and test data

-Find best Hyperparameter for SVM, Classification Trees and Logistic Regression

- Find the method performs best using test data

# Import Libraries and Define Auxiliary Functions

We will import the following libraries for the lab

```
In [3]:
%%capture
# Pandas is a software library written for the Python programming language for
data manipulation and analysis.
import pandas as pd
# NumPy is a library for the Python programming language, adding support for
large, multi-dimensional arrays and matrices, along with a large collection of
high-level mathematical functions to operate on these arrays
import numpy as np
# Matplotlib is a plotting library for python and pyplot gives us a MatLab like
plotting framework. We will use this in our plotter function to plot data.
import matplotlib.pyplot as plt
#Seaborn is a Python data visualization library based on matplotlib. It provides
a high-level interface for drawing attractive and informative statistical
graphics
import seaborn as sns

!pip install -U scikit-learn
```

```
!pip install -U js
!pip install --upgrade scikit-learn

# Preprocessing allows us to standarsize our data
from sklearn import preprocessing
# Allows us to split our data into training and testing data
from sklearn.model_selection import train_test_split
# Allows us to test parameters of classification algorithms and find the best one
from sklearn.model_selection import GridSearchCV
# Logistic Regression classification algorithm
from sklearn.linear_model import LogisticRegression
# Support Vector Machine classification algorithm
from sklearn.svm import SVC
# Decision Tree classification algorithm
from sklearn.tree import DecisionTreeClassifier
# K Nearest Neighbors classification algorithm
from sklearn.neighbors import KNeighborsClassifier
```

This function is to plot the confusion matrix.

```
In [4]:  def plot_confusion_matrix(y,y_predict):
             "this function plots the confusion matrix"
             from sklearn.metrics import confusion_matrix

             cm = confusion_matrix(y, y_predict)
             ax= plt.subplot()
             sns.heatmap(cm, annot=True, ax = ax); #annot=True to annotate cells
             ax.set_xlabel('Predicted labels')
             ax.set_ylabel('True labels')
             ax.set_title('Confusion Matrix');
             ax.xaxis.set_ticklabels(['did not land', 'land']);
         ax.yaxis.set_ticklabels(['did not land', 'landed'])
             plt.show()
```

# Load the dataframe

Load the data

```
In [5]:  import requests
         import io
         import pandas as pd

         URL1 = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-
         DS0321EN-SkillsNetwork/datasets/dataset_part_2.csv"
         resp1 = requests.get(URL1)
         text1 = io.StringIO(resp1.text)
         data = pd.read_csv(text1)
```

```
In [6]:  data.head()
```

Out[6]:

| | FlightNumber | Date | BoosterVersion | PayloadMass | Orbit | LaunchSite | Outcome | Flight |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 2010-06-04 | Falcon 9 | 6104.959412 | LEO | CCAFS SLC 40 | None None | . |
| **1** | 2 | 2012-05-22 | Falcon 9 | 525.000000 | LEO | CCAFS SLC 40 | None None | . |
| **2** | 3 | 2013-03-01 | Falcon 9 | 677.000000 | ISS | CCAFS SLC 40 | None None | . |
| **3** | 4 | 2013-09-29 | Falcon 9 | 500.000000 | PO | VAFB SLC 4E | False Ocean | . |
| **4** | 5 | 2013-12-03 | Falcon 9 | 3170.000000 | GTO | CCAFS SLC 40 | None None | . |

In [7]:
```
URL2 = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-
DS0321EN-SkillsNetwork/datasets/dataset_part_3.csv'
resp2 = requests.get(URL2)
text2 = io.BytesIO(resp2.content)
X = pd.read_csv(text2)
```

In [8]:
```
X.head(100)
```

Out[8]:

| | FlightNumber | PayloadMass | Flights | Block | ReusedCount | Orbit_ES-L1 | Orbit_GEO | Orbit |
|---|---|---|---|---|---|---|---|---|
| **0** | 1.0 | 6104.959412 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | |
| **1** | 2.0 | 525.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | |
| **2** | 3.0 | 677.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | |
| **3** | 4.0 | 500.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | |
| **4** | 5.0 | 3170.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | |
| **85** | 86.0 | 15400.000000 | 2.0 | 5.0 | 2.0 | 0.0 | 0.0 | |
| **86** | 87.0 | 15400.000000 | 3.0 | 5.0 | 2.0 | 0.0 | 0.0 | |
| **87** | 88.0 | 15400.000000 | 6.0 | 5.0 | 5.0 | 0.0 | 0.0 | |
| **88** | 89.0 | 15400.000000 | 3.0 | 5.0 | 2.0 | 0.0 | 0.0 | |
| **89** | 90.0 | 3681.000000 | 1.0 | 5.0 | 0.0 | 0.0 | 0.0 | |

90 rows × 83 columns

# TASK 1

Create a NumPy array from the column `Class` in `data`, by applying the method `to_numpy()` then assign it to the variable `Y`, make sure the output is a Pandas series (only one bracket df['name of column']).

```python
In [9]:  Y = data['Class'].to_numpy()
```

## TASK 2

Standardize the data in `X` then reassign it to the variable `X` using the transform provided below.

```python
In [10]:  X.mean()
          X.std()
```

```
Out[10]:  FlightNumber          26.124701
          PayloadMass         4694.671720
          Flights                1.213172
          Block                  1.595288
          ReusedCount            1.710254
                                ...
          GridFins_True          0.418069
          Reused_False           0.494792
          Reused_True            0.494792
          Legs_False             0.410383
          Legs_True              0.410383
          Length: 83, dtype: float64
```

```python
In [11]:  # students get this
          transform = preprocessing.StandardScaler().fit(X)
```

We split the data into training and testing data using the function `train_test_split`. The training data is divided into validation data, a second set used for training data; then the models are trained and hyperparameters are selected using the function `GridSearchCV`.

## TASK 3

Use the function train_test_split to split the data X and Y into training and test data. Set the parameter test_size to 0.2 and random_state to 2. The training data and test data should be assigned to the following labels.

`X_train, X_test, Y_train, Y_test`

```python
In [12]:  X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2,
          random_state=2)
```

we can see we only have 18 test samples.

```
In [13]:  Y_test.shape
```

```
Out[13]:  (18,)
```

## TASK 4

Create a logistic regression object then create a GridSearchCV object `logreg_cv` with cv = 10. Fit the object to find the best parameters from the dictionary `parameters`.

```
In [14]:  parameters ={'C':[0.01,0.1,1],
                       'penalty':['l2'],
                       'solver':['lbfgs']}
```

```
In [15]:  parameters ={"C":[0.01,0.1,1],'penalty':['l2'], 'solver':['lbfgs']}# l1 lasso l2
          ridge
          lr=LogisticRegression()
```

We output the `GridSearchCV` object for logistic regression. We display the best parameters using the data attribute `best_params_` and the accuracy on the validation data using the data attribute `best_score_`.

```
In [16]:  %%capture
          logreg_cv = GridSearchCV(lr, parameters, cv=10)
          logreg_cv.fit(X_train, Y_train)
```

```
In [17]:  print("tuned hpyerparameters :(best parameters) ",logreg_cv.best_params_)
          print("accuracy :",logreg_cv.best_score_)

          tuned hpyerparameters :(best parameters)  {'C': 1, 'penalty': 'l2', 'solver': 'lbfg
          s'}
          accuracy : 0.8196428571428571
```
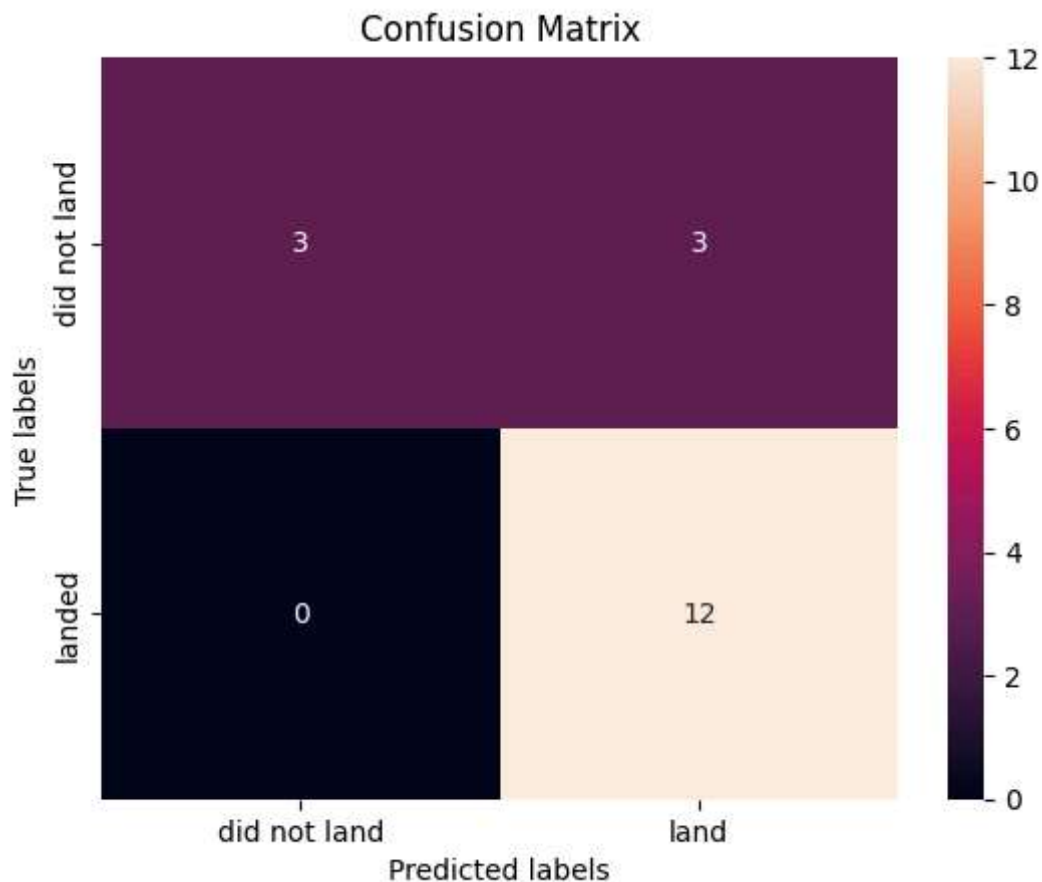
## TASK 5

Calculate the accuracy on the test data using the method `score`:

```
In [18]:  logreg_cv.score(X_test, Y_test)
```

```
Out[18]:  0.8333333333333334
```

Lets look at the confusion matrix:

```
In [19]:  yhat=logreg_cv.predict(X_test)
          plot_confusion_matrix(Y_test,yhat)
```

Confusion Matrix

Examining the confusion matrix, we see that logistic regression can distinguish between the different classes. We see that the major problem is false positives.

## TASK 6

Create a support vector machine object then create a `GridSearchCV` object `svm_cv` with cv - 10. Fit the object to find the best parameters from the dictionary `parameters`.

```
In [ ]: parameters = {'kernel':('linear', 'rbf','poly', 'sigmoid'),
                      'C': np.logspace(-3, 3, 5),
                      'gamma':np.logspace(-3, 3, 5)}
        svm = SVC()
```

```
In [ ]: svm_cv = GridSearchCV(svm, parameters, cv=3)#10)
        svm_cv.fit(X_train, Y_train)
```

```
In [ ]: print("tuned hpyerparameters :(best parameters) ",svm_cv.best_params_)
        print("accuracy :",svm_cv.best_score_)
```

## TASK 7

Calculate the accuracy on the test data using the method `score` :

```
In [ ]:  svm_cv.score(X_test, Y_test)
```

We can plot the confusion matrix

```
In [ ]:  yhat=svm_cv.predict(X_test)
         plot_confusion_matrix(Y_test,yhat)
```

## TASK 8

Create a decision tree classifier object then create a `GridSearchCV` object `tree_cv` with cv = 10. Fit the object to find the best parameters from the dictionary `parameters`.

```
In [ ]:  parameters = {'criterion': ['gini', 'entropy'],
              'splitter': ['best', 'random'],
              'max_depth': [2*n for n in range(1,10)],
              'max_features': ['auto', 'sqrt'],
              'min_samples_leaf': [1, 2, 4],
              'min_samples_split': [2, 5, 10]}

         tree = DecisionTreeClassifier()
```

```
In [ ]:  tree_cv = GridSearchCV(tree, parameters, cv=10)
```

```
In [ ]:  print("tuned hpyerparameters :(best parameters) ",tree_cv.best_params_)
         print("accuracy :",tree_cv.best_score_)
```

## TASK 9

Calculate the accuracy of tree_cv on the test data using the method `score`:

```
In [ ]:  tree_cv.score(X_test,Y_test)
```

We can plot the confusion matrix

```
In [ ]:  yhat = tree_cv.predict(X_test)
         plot_confusion_matrix(Y_test,yhat)
```

## TASK 10

Create a k nearest neighbors object then create a `GridSearchCV` object `knn_cv` with cv = 10. Fit the object to find the best parameters from the dictionary `parameters`.

```
In [ ]:  parameters = {'n_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
                  'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute'],
                  'p': [1,2]}
```

```
KNN = KNeighborsClassifier()
```

In [ ]:
```
knn_cv = GridSearchCV(KNN, parameters, cv=10)
knn_cv.fit(X_train,Y_train)
```

In [ ]:
```
print("tuned hpyerparameters :(best parameters) ",knn_cv.best_params_)
print("accuracy :",knn_cv.best_score_)
```

## TASK 11

Calculate the accuracy of knn_cv on the test data using the method `score` :

In [ ]:
```
knn_cv.score(X_test,Y_test)
```

We can plot the confusion matrix

In [ ]:
```
yhat = knn_cv.predict(X_test)
plot_confusion_matrix(Y_test,yhat)
```

## TASK 12

Find the method performs best:

Similar results from all algorithms.

## Authors

Pratiksha Verma

## Change Log

| Date (YYYY-MM-DD) | Version | Changed By | Change Description |
|---|---|---|---|
| 2022-11-09 | 1.0 | Pratiksha Verma | Converted initial version to Jupyterlite |