# 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:



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 Librairies and Define Auxiliary Functions

```
[1]: import piplite
     await piplite.install(['numpy'])
     await piplite.install(['pandas'])
     await piplite.install(['seaborn'])
```

We will import the following libraries for the lab

```
[2]: # 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, al
     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 p
     import matplotlib.pyplot as plt
     #Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractiv
     import seaborn as sns
     # 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.

```python
[3]: 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

```python
[4]: from js import fetch
     import io

     URL1 = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_2
     resp1 = await fetch(URL1)
     text1 = io.BytesIO((await resp1.arrayBuffer()).to_py())
     data = pd.read_csv(text1)
```

```
[5]: data.head()
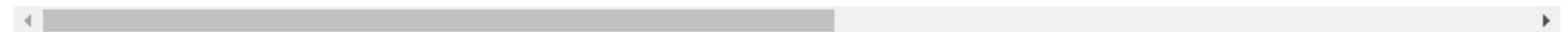```

| | FlightNumber | Date | BoosterVersion | PayloadMass | Orbit | LaunchSite | Outcome | Flights | GridFins | Reused | Legs | LandingPad | Block |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 2010-06-04 | Falcon 9 | 6104.959412 | LEO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 |
| **1** | 2 | 2012-05-22 | Falcon 9 | 525.000000 | LEO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 |
| **2** | 3 | 2013-03-01 | Falcon 9 | 677.000000 | ISS | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 |
| **3** | 4 | 2013-09-29 | Falcon 9 | 500.000000 | PO | VAFB SLC 4E | False Ocean | 1 | False | False | False | NaN | 1.0 |
| **4** | 5 | 2013-12-03 | Falcon 9 | 3170.000000 | GTO | CCAFS SLC 40 | None None | 1 | False | False | False | NaN | 1.0 |

```
[6]: URL2 = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_3
resp2 = await fetch(URL2)
text2 = io.BytesIO((await resp2.arrayBuffer()).to_py())
X = pd.read_csv(text2)
```

```
[7]: X.head(100)
```

[7]:

| | FlightNumber | PayloadMass | Flights | Block | ReusedCount | Orbit_ES-L1 | Orbit_GEO | Orbit_GTO | Orbit_HEO | Orbit_ISS | ... | Serial_B1058 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.0 | 6104.959412 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 |
| 1 | 2.0 | 525.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 |
| 2 | 3.0 | 677.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | ... | 0.0 |
| 3 | 4.0 | 500.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 |
| 4 | 5.0 | 3170.000000 | 1.0 | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 85 | 86.0 | 15400.000000 | 2.0 | 5.0 | 2.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 |
| 86 | 87.0 | 15400.000000 | 3.0 | 5.0 | 2.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 1.0 |
| 87 | 88.0 | 15400.000000 | 6.0 | 5.0 | 5.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 |
| 88 | 89.0 | 15400.000000 | 3.0 | 5.0 | 2.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 |
| 89 | 90.0 | 3681.000000 | 1.0 | 5.0 | 0.0 | 0.0 | 0.0 | 0.0 | 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']).

```
[8]: Y = data["Class"].to_numpy()
```

# TASK 2

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

```
[9]: # students get this
transform = preprocessing.StandardScaler()
X = transform.fit_transform(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`

```
[10]: 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.

```
[11]: Y_test.shape
```

```
[11]: (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`.

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

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

```
[14]: gsvc=GridSearchCV(lr, parameters, scoring = 'accuracy',cv=10)
     logreg_cv =gsvc.fit(X_train, Y_train)
```

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_`.

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

```
tuned hpyerparameters :(best parameters)  {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
accuracy : 0.8464285714285713
```
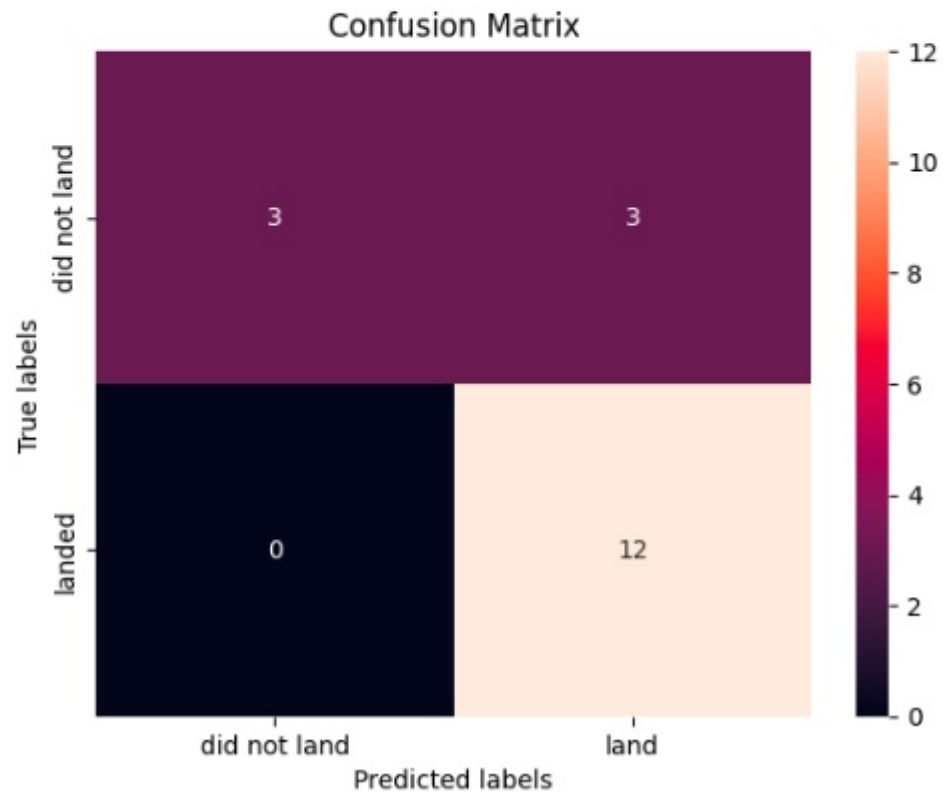
# TASK 5

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

```
[16]: print('Accuracy = ', logreg_cv.score(X_test, Y_test))
```

```
Accuracy =  0.8333333333333334
```

Lets look at the confusion matrix:

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

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` .

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

```
[20]: gscv = GridSearchCV(svm,parameters, scoring= 'accuracy', cv=10)
      svm_cv = gscv.fit(X_train, Y_train)
```

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

```
tuned hpyerparameters :(best parameters)  {'C': 1.0, 'gamma': 0.03162277660168379, 'kernel': 'sigmoid'}
accuracy : 0.8482142857142856
```
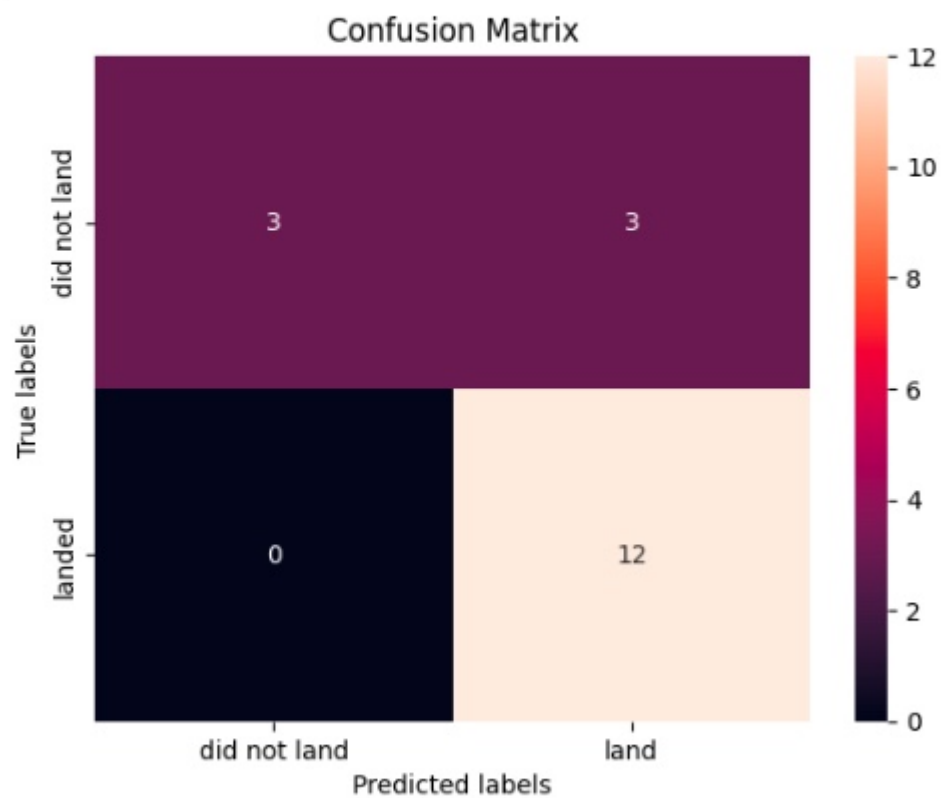
## TASK 7

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

```
[22]: print('accuracy: ', svm_cv.score(X_test, Y_test))
```

```
accuracy:  0.8333333333333334
```

We can plot the confusion matrix

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

## Confusion Matrix

# 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`.

```python
[24]: 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()
```

```python
[26]: gscv = GridSearchCV(tree,parameters, scoring='accuracy',cv=10)
      tree_cv = gsvc.fit(X_train, Y_train)
```

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

```
---------------------------------------------------------------------------
NameError                                 Traceback (most recent call last)
Cell In[24], line 1
----> 1 print("tuned hpyerparameters :(best parameters) ",tree_cv.best_params_)
      2 print("accuracy :",tree_cv.best_score_)

NameError: name 'tree_cv' is not defined
```
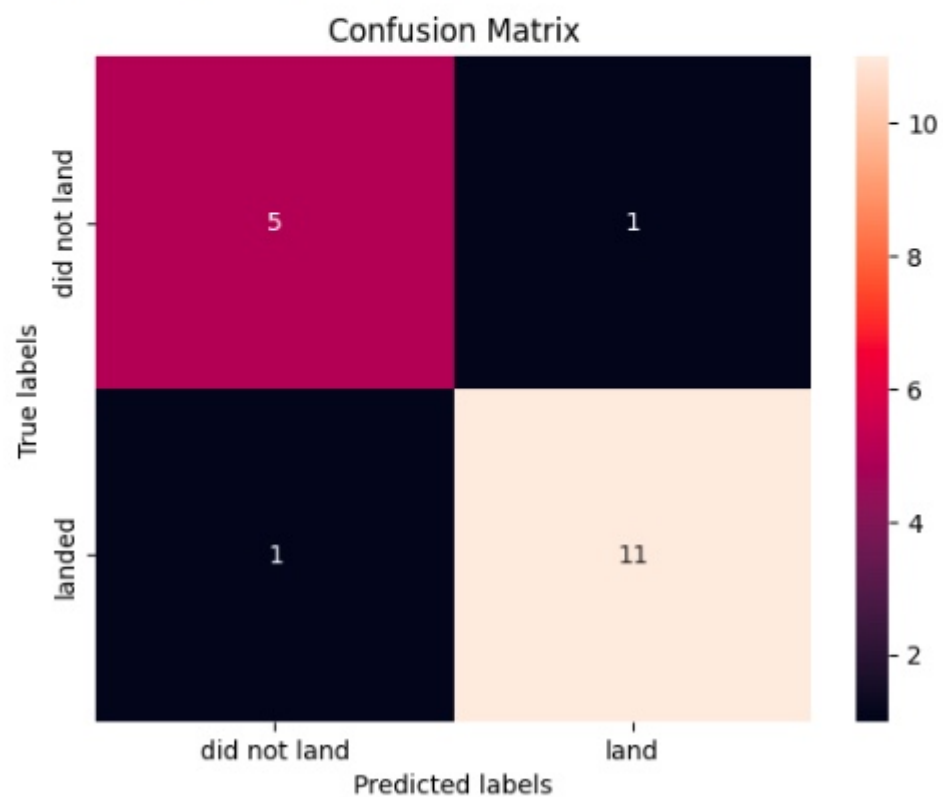
# TASK 9

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

```python
[29]: print("accuracy: ", tree_cv.score(X_test, Y_test))

accuracy:  0.8888888888888888
```

We can plot the confusion matrix

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

### Confusion Matrix

## 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`.

```
[31]: 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()
```

```
[32]: gscv = GridSearchCV(KNN,parameters,scoring='accuracy',cv=10)
      knn_cv = gscv.fit(X_train, Y_train)
```

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

```
tuned hpyerparameters :(best parameters)  {'algorithm': 'auto', 'n_neighbors': 10, 'p': 1}
accuracy : 0.8482142857142858
```

## TASK 11
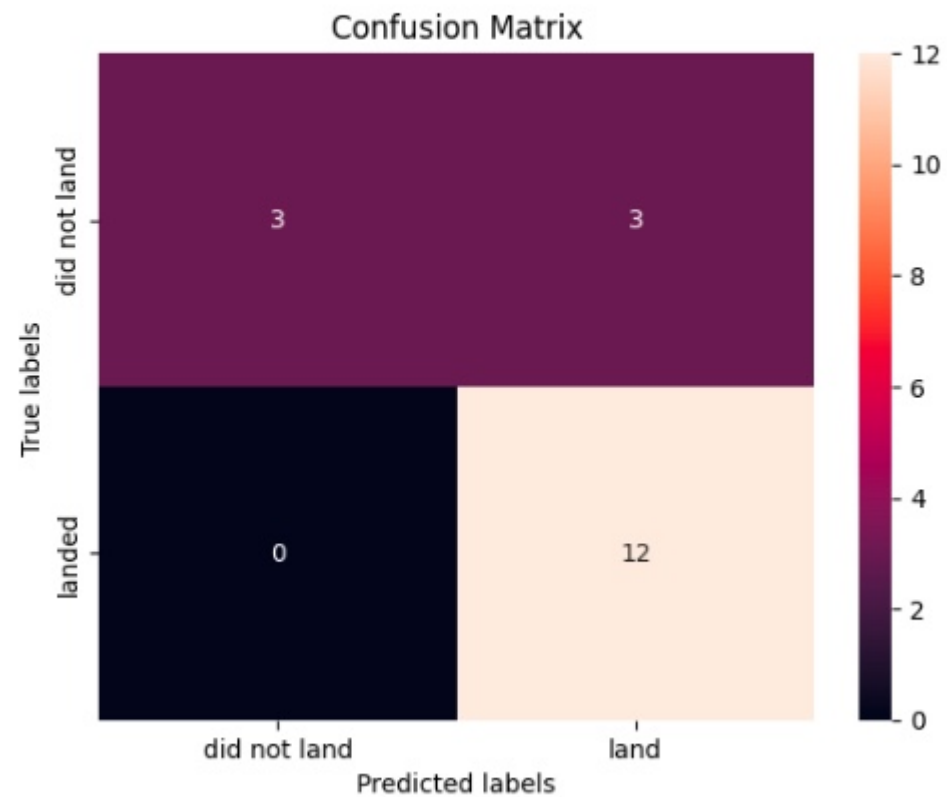
Calculate the accuracy of knn_cv on the test data using the method `score`:

```
[34]: print("accuracy: ", knn_cv.score(X_test, Y_test))
```

```
accuracy:  0.8333333333333334
```

We can plot the confusion matrix

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

## TASK 12

Find the method performs best:

```python
[36]: algorithms = {'KNN':knn_cv.best_score_,'Tree':tree_cv.best_score_,'LogisticRegression':logreg_cv.best_score_}
      bestalgorithm = max(algorithms, key=algorithms.get)
      print('Best Algorithm is',bestalgorithm,'with a score of',algorithms[bestalgorithm])
      if bestalgorithm == 'KNN':
          print('Best Params is :',knn_cv.best_params_)
      if bestalgorithm == 'LogisticRegression':
          print('Best Params is :',logreg_cv.best_params_)
```

Best Algorithm is Tree with a score of 0.8642857142857142

## 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 |