

## 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.

Requirement already satisfied: numpy>=1.17.3; platform\_machine != "aarch64" and platform\_machine != "arm64" and python\_version < "3.10" in /opt/conda/lib/python3.7/site-packages (from panda

# 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

Outcome Flights GridFins Reused Legs LandingPad Block ReusedCount Serial

NaN

NaN

NaN

NaN

NaN

0.0

0.0

0.0

0.0

0.0

0.0

1.0

0.0

0.0

0.0

1.0

1.0

1.0

1.0

1.0

Orbit\_GEO Orbit\_GTO Orbit\_HEO Orbit\_ISS ... Serial\_B1058 Serial\_B1059 Serial\_B1060 Serial\_B1062 GridFins\_False GridFins\_True Reused\_False Reu

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

False

False

False

False

False

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

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

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.

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\_.

False

False False

False False

False False

False False

0.0 ...

0.0 ...

1.0 ...

0.0 ...

0.0 ...

0.0 ...

0.0 ...

0.0 ...

0.0 ...

0.0 ...

Latitude Class

1.0

1.0

1.0

1.0

1.0

0.0

0.0

0.0

0.0

0.0

0

0

0

0.0

0.0

0.0

0.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

0.0

0.0

0.0

0.0

1.0

Longitude

0 B0003

0 B0005

0 B0007

0 B1003

0 B1004

0.0

0.0

0.0

1.0

0.0

0.0

1.0

0.0

-80.577366 28.561857

-80.577366 28.561857

-80.577366 28.561857

-120.610829 34.632093

-80.577366 28.561857

0.0

0.0

0.0

0.0

0.0

0.0

0.0

1.0

# 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.

#Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics

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

In [5]: !pip install numpy !pip install pandas !pip install seaborn Requirement already satisfied: numpy in /opt/conda/lib/python3.7/site-packages (1.21.6) Requirement already satisfied: pandas in /opt/conda/lib/python3.7/site-packages (1.3.5)

s) (1.21.6) Requirement already satisfied: python-dateutil>=2.7.3 in /opt/conda/lib/python3.7/site-packages (from pandas) (2.8.1) Requirement already satisfied: pytz>=2017.3 in /opt/conda/lib/python3.7/site-packages (from pandas) (2024.2) Requirement already satisfied: six>=1.5 in /opt/conda/lib/python3.7/site-packages (from python-dateutil>=2.7.3->pandas) (1.15.0) Requirement already satisfied: seaborn in /opt/conda/lib/python3.7/site-packages (0.12.2)

Requirement already satisfied: typing\_extensions; python\_version < "3.8" in /opt/conda/lib/python3.7/site-packages (from seaborn) (4.7.1)

Requirement already satisfied: matplotlib!=3.6.1,>=3.1 in /opt/conda/lib/python3.7/site-packages (from seaborn) (3.5.3) Requirement already satisfied: numpy!=1.24.0,>=1.17 in /opt/conda/lib/python3.7/site-packages (from seaborn) (1.21.6) Requirement already satisfied: pandas>=0.25 in /opt/conda/lib/python3.7/site-packages (from seaborn) (1.3.5) Requirement already satisfied: cycler>=0.10 in /opt/conda/lib/python3.7/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (0.11.0) Requirement already satisfied: python-dateutil>=2.7 in /opt/conda/lib/python3.7/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (2.8.1) Requirement already satisfied: packaging>=20.0 in /opt/conda/lib/python3.7/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (20.4)

Requirement already satisfied: kiwisolver>=1.0.1 in /opt/conda/lib/python3.7/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (1.4.5) Requirement already satisfied: pillow>=6.2.0 in /opt/conda/lib/python3.7/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (9.5.0) Requirement already satisfied: fonttools>=4.22.0 in /opt/conda/lib/python3.7/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (4.38.0) Requirement already satisfied: pyparsing>=2.2.1 in /opt/conda/lib/python3.7/site-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (2.4.7)

Requirement already satisfied: pytz>=2017.3 in /opt/conda/lib/python3.7/site-packages (from pandas>=0.25->seaborn) (2024.2) Requirement already satisfied: six>=1.5 in /opt/conda/lib/python3.7/site-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.1->seaborn) (1.15.0) We will import the following libraries for the lab

In [7]: !pip install scikit-learn Collecting scikit-learn Downloading scikit\_learn-1.0.2-cp37-cp37m-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl (24.8 MB) | 24.8 MB 33.6 MB/s eta 0:00:01

Collecting threadpoolctl>=2.0.0

Collecting joblib>=0.11

import pandas as pd

import numpy as np

import seaborn as sns

import matplotlib.pyplot as plt

from sklearn.svm import SVC

| 302 kB 82.6 MB/s eta 0:00:01 Requirement already satisfied: numpy>=1.14.6 in /opt/conda/lib/python3.7/site-packages (from scikit-learn) (1.21.6) Collecting scipy>=1.1.0 Downloading scipy-1.7.3-cp37-cp37m-manylinux\_2\_12\_x86\_64.manylinux2010\_x86\_64.whl (38.1 MB) | 38.1 MB 87.1 MB/s eta 0:00:01 Installing collected packages: threadpoolctl, joblib, scipy, scikit-learn Successfully installed joblib-1.3.2 scikit-learn-1.0.2 scipy-1.7.3 threadpoolctl-3.1.0

Downloading threadpoolctl-3.1.0-py3-none-any.whl (14 kB)

In [8]: # Pandas is a software library written for the Python programming language for data manipulation and analysis.

ax.xaxis.set\_ticklabels(['did not land', 'land']); ax.yaxis.set\_ticklabels(['did not land', 'landed'])

URL = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset\_part\_2.csv"

LaunchSite

PO VAFB SLC 4E False Ocean

None None

None None

6104.959412 LEO CCAFS SLC 40

525.000000 LEO CCAFS SLC 40

Falcon 9 3170.000000 GTO CCAFS SLC 40 None None

ISS CCAFS SLC 40

In [14]: URL2 = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset\_part\_3.csv'

L1

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

2.0

2.0

5.0

2.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

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.

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

param\_grid={'C': array([1.00000000e-03, 3.16227766e-02, 1.00000000e+00, 3.16227766e+01,

'kernel': ('linear', 'rbf', 'poly', 'rbf', 'sigmoid')})

'gamma': array([1.00000000e-03, 3.16227766e-02, 1.00000000e+00, 3.16227766e+01,

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.

0.0

0.0

0.0

1.0

0.0

0.0

0.0

0.0

0.0

Downloading joblib-1.3.2-py3-none-any.whl (302 kB)

# 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.tree import DecisionTreeClassifier # K Nearest Neighbors classification algorithm from sklearn.neighbors import KNeighborsClassifier

# Decision Tree classification algorithm

In [9]: 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');

plt.show()

Load the data

data.head(5)

FlightNumber

In [13]: **import** requests import io

Load the dataframe

response1 = requests.get(URL)

data = pd.read\_csv(text1)

text1 = io.BytesIO(response1.content)

Date BoosterVersion PayloadMass Orbit

677.000000

500.000000

Falcon 9

Falcon 9

Falcon 9

1.0

1.0

1.0

5.0

5.0

5.0

5.0

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

In [21]: X\_train, X\_test, Y\_train, Y\_test=train\_test\_split(X, Y, test\_size=0.2, random\_state=2)

In [26]: parameters ={"C":[0.01,0.1,1],'penalty':['l2'], 'solver':['lbfgs']}# 11 lasso 12 ridge

param\_grid={'C': [0.01, 0.1, 1], 'penalty': ['12'],

'solver': ['lbfgs']})

In [29]: print("Logistic Regression test data accuracy :",logreg\_cv.score(X\_test, Y\_test))

1.0

1.0

1.0

1.0

2.0

3.0

6.0

3.0

1.0

500.000000

5.0 3170.000000

86.0 15400.000000

87.0 15400.000000

88.0 15400.000000

89.0 15400.000000

90.0 3681.000000

transform = preprocessing.StandardScaler()

selected using the function GridSearchCV.

X\_train, X\_test, Y\_train, Y\_test

we can see we only have 18 test samples.

'penalty':['12'], 'solver':['lbfgs']}

# Create a GridSearchCV object logreg\_cv

Out[26]: GridSearchCV(cv=10, estimator=LogisticRegression(),

logreg\_cv = GridSearchCV(lr, parameters, cv=10)

#Fit the training data into the GridSearch object

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

Logistic Regression test data accuracy : 0.833333333333333333

This function is to plot the confusion matrix.

0 1 2010-06-04 2 2012-05-22 2 3 2013-03-01 4 2013-09-29 4 5 2013-12-03

resp2 = requests.get(URL2)

X = pd.read\_csv(text2)

In [15]: X

text2 = io.BytesIO(resp2.content)

FlightNumber PayloadMass Flights Block ReusedCount 1.0 6104.959412 525.000000 2 677.000000

3

85

86

87

88

89

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']). In [20]: Y = data['Class'].to\_numpy()

In [18]: X = transform.fit(X).transform(X)

90 rows × 83 columns

TASK 2

In [17]: # students get this

TASK 3

In [22]: Y\_test.shape Out [22]: (18,) TASK 4

In [23]: parameters ={'C':[0.01,0.1,1],

lr=LogisticRegression()

logreg\_cv.fit(X\_train, Y\_train)

Lets look at the confusion matrix:

plot\_confusion\_matrix(Y\_test,yhat)

Confusion Matrix

Predicted labels

In [32]: # create a GridSearchCV object svm\_cv with cv = 10 svm\_cv = GridSearchCV(svm, parameters, cv= 10)

svm\_cv.fit(X\_train, Y\_train)

Out[32]: GridSearchCV(cv=10, estimator=SVC(),

#Fit the training data into the GridSearch object

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

In [34]: print("SVM test data accuracy :",svm\_cv.score(X\_test, Y\_test))

12

land

- 12

SVM test data accuracy : 0.8333333333333333

Confusion Matrix

Predicted labels

In [36]: parameters = {'criterion': ['gini', 'entropy'], 'splitter': ['best', 'random'],

tree = DecisionTreeClassifier()

tree\_cv.fit(X\_train, Y\_train)

accuracy : 0.8892857142857142

We can plot the confusion matrix

plot\_confusion\_matrix(Y\_test,yhat)

Confusion Matrix

Predicted labels

'p': [1,2]}

knn\_cv.fit(X\_train, Y\_train)

#Fit the training data into the GridSearch object

'p': [1, 2]})

In [45]: print("tuned hpyerparameters : (best parameters) ",knn\_cv.best\_params\_)

Out[44]: GridSearchCV(cv=10, estimator=KNeighborsClassifier(),

print("accuracy :",knn\_cv.best\_score\_)

In [42]: yhat = tree\_cv.predict(X\_test)

did not land

'max\_features': ['auto', 'sqrt'], 'min\_samples\_leaf': [1, 2, 4], 'min\_samples\_split': [2, 5, 10]}

In [39]: # create a GridSearchCV object tree\_cv with cv = 10

tree\_cv = GridSearchCV(tree, parameters, cv= 10)

#Fit the training data into the GridSearch object

param\_grid={'criterion': ['gini', 'entropy'],

In [40]: print("tuned hpyerparameters : (best parameters) ",tree\_cv.best\_params\_)

'max\_features': ['auto', 'sqrt'], 'min\_samples\_leaf': [1, 2, 4], 'min\_samples\_split': [2, 5, 10], 'splitter': ['best', 'random']})

'algorithm': ['auto', 'ball\_tree', 'kd\_tree', 'brute'],

param\_grid={'algorithm': ['auto', 'ball\_tree', 'kd\_tree', 'brute'], 'n\_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],

tuned hpyerparameters : (best parameters) {'algorithm': 'auto', 'n\_neighbors': 10, 'p': 1}

'max\_depth': [2, 4, 6, 8, 10, 12, 14, 16, 18],

tuned hpyerparameters: (best parameters) {'criterion': 'entropy', 'max\_depth': 18, 'max\_features': 'sqrt', 'min\_samples\_leaf': 4, 'min\_samples\_split': 2, 'splitter': 'random'}

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.

Out[39]: GridSearchCV(cv=10, estimator=DecisionTreeClassifier(),

print("accuracy :", tree\_cv.best\_score\_)

'max\_depth': [2\*n for n in range(1,10)],

We can plot the confusion matrix

plot\_confusion\_matrix(Y\_test,yhat)

In [35]: yhat=svm\_cv.predict(X\_test)

In [31]: parameters = {'kernel':('linear', 'rbf', 'poly', 'rbf', 'sigmoid'), 'C': np.logspace(-3, 3, 5), 'gamma':np.logspace(-3, 3, 5)}

In [30]: yhat=logreg\_cv.predict(X\_test)

did not land

In [27]: print("tuned hpyerparameters : (best parameters) ",logreg\_cv.best\_params\_) print("accuracy :", logreg\_cv.best\_score\_) tuned hpyerparameters :(best parameters) {'C': 0.01, 'penalty': '12', 'solver': 'lbfgs'} accuracy : 0.8464285714285713 TASK 5

TASK 6

svm = SVC()

1.00000000e+03]), In [33]: 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

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.

did not land

TASK 9 Calculate the accuracy of tree\_cv on the test data using the method score: In [41]: print("Decision Tree accuracy on test set :",tree\_cv.score(X\_test, Y\_test)) 

In [43]: parameters = {'n\_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10], KNN = KNeighborsClassifier() In [44]: # create a GridSearchCV object knn\_cv with cv = 10 knn\_cv = GridSearchCV(KNN, parameters, cv= 10)

TASK 10

accuracy : 0.8482142857142858 TASK 11 Calculate the accuracy of knn\_cv on the test data using the method score: In [46]: print("K nearest neighbors accuracy on test set :",knn\_cv.score(X\_test, Y\_test))

Logistic\_Reg

**Decision Tree** 

**Authors** 

Pratiksha Verma

<!--## Change Log--!>

SVM

KNN

Out[48]:

did not land Predicted labels **TASK 12** Find the method performs best: Report.transpose()

12 land In [48]: Report = pd.DataFrame({'Method' : ['Test Data Accuracy']}) knn\_accuracy=knn\_cv.score(X\_test, Y\_test) Decision\_tree\_accuracy=tree\_cv.score(X\_test, Y\_test) SVM\_accuracy=svm\_cv.score(X\_test, Y\_test) Logistic\_Regression=logreg\_cv.score(X\_test, Y\_test) Report['Logistic\_Reg'] = [Logistic\_Regression] Report['SVM'] = [SVM\_accuracy] Report['Decision Tree'] = [Decision\_tree\_accuracy] Report['KNN'] = [knn\_accuracy]

Method Test Data Accuracy

K nearest neighbors accuracy on test set : 0.83333333333333333 We can plot the confusion matrix

In [47]: yhat = knn\_cv.predict(X\_test) plot\_confusion\_matrix(Y\_test,yhat) Confusion Matrix

- 12

0

0.833333

0.833333

0.888889

0.833333

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

IBM Corporation 2022. All rights reserved.