

Assignment: Machine Learning Prediction

Several examples of an unsuccessful landing are shown here:

Perform exploratory Data Analysis and determine Training Labels

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

Import Libraries and Define Auxiliary Functions

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

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

SEPTEMBER 2013 HARD IMPACT ON OCEAN

create a column for the class

Split into training data and test data

await piplite.install(['numpy']) await piplite.install(['pandas']) await piplite.install(['seaborn'])

import matplotlib.pyplot as plt

from sklearn import preprocessing

from sklearn.svm import SVC

import pandas as pd

import numpy as np

import seaborn as sns

We will import the following libraries for the lab

Preprocessing allows us to standarsize our data

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

"this function plots the confusion matrix" from sklearn.metrics import confusion_matrix

text1 = io.BytesIO((await resp1.arrayBuffer()).to_py())

Date BoosterVersion PayloadMass Orbit

677.000000

500.000000

Falcon 9

Falcon 9

Falcon 9

Falcon 9

PayloadMass Flights Block ReusedCount

1.0

1.0

1.0

1.0

1.0

2.0

6.0

3.0

1.0

3.0 5.0

1.0

1.0

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

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

print("tuned hpyerparameters :(best parameters) ",logreg_cv.best_params_)

In []: from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

metrics = ['log_accuracy', 'log_precision', 'log_recall', 'log_f1']

df_sorted2 = df_sorted.sort_values(by = 'score', ascending=False)

score = [log_accuracy,log_precision, log_recall, log_f1]

sns.barplot(x = 'metrics', y = 'score', data = df_sorted2)

classification metrics

Metrics

print(log_accuracy,log_precision, log_recall, log_f1)

12

land

parameters = {'kernel':('linear', 'rbf', 'poly', 'rbf', 'sigmoid'),

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

SVC

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

metrics = ['svm_accuracy', 'svm_precision', 'svm_recall', 'svm_f1']

df_sorted_svm = pd.DataFrame(data).sort_values(by='score', ascending=False)

sns.barplot(x='metrics', y='score', data=df_sorted_svm, palette='magma')

plt.xticks(rotation=45) # Rotate x-axis labels for better visibility

sns.barplot(x='metrics', y='score', data=df_sorted_svm, palette='magma')

SVM Classification Metrics

Metrics

- 12

score = [svm_accuracy, svm_precision, svm_recall, svm_f1]

plt.title('SVM Classification Metrics', fontsize=16)

svm_accuracy = accuracy_score(Y_test, Y_hat_svm) svm_precision = precision_score(Y_test, Y_hat_svm)

svm_recall = recall_score(Y_test, Y_hat_svm)

Store metrics and scores in a DataFrame

data = {'metrics': metrics, 'score': score}

svm_f1 = f1_score(Y_test, Y_hat_svm)

SVC(gamma=0.03162277660168379, kernel='sigmoid')

tuned hpyerparameters :(best parameters) {'C': 1.0, 'gamma': 0.03162277660168379, 'kernel': 'sigmoid'}

'C': np.logspace(-3, 3, 5), 'gamma':np.logspace(-3, 3, 5)}

svm_cv = GridSearchCV(svm, parameters, cv=10)

print("accuracy :", svm_cv.best_score_)

accuracy: 0.8482142857142856

In []: Y_hat_svm = svm_cv.predict(X_test) # Calculate evaluation metrics

Plot the bar plot

plt.show()

0.8

0.4

0.2

True labels

We can plot the confusion matrix

did not land

TASK 8

In []: plot_confusion_matrix(Y_test,Y_hat_svm)

Confusion Matrix

Predicted labels

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

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

print("accuracy :", tree_cv.best_score_)

score = accuracy_score(Y_test, Y_hat)

In []: from sklearn.metrics import accuracy_score

plot_confusion_matrix(Y_test, yhat)

Confusion Matrix

Predicted labels

'p': [1,2]}

GridSearchCV

▶ estimator: KNeighborsClassifier

KNeighborsClassifier

print("accuracy :", knn_cv.best_score_)

knn_cv = GridSearchCV(KNN, parameters, cv=10)

KNN = KNeighborsClassifier()

knn_cv.fit(X_train, Y_train)

accuracy : 0.8482142857142858

knn_cv.score(X_test, Y_test)

We can plot the confusion matrix

plot_confusion_matrix(Y_test,yhat)

Confusion Matrix

Predicted labels

In []: yhat = knn_cv.predict(X_test)

did not land

Find the method performs best:

TASK 12

Authors

Pratiksha Verma

Change Log

TASK 11

Out[]: 0.8333333333333333

did not land

True labels

12

In []: parameters = {'n_neighbors': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],

We can plot the confusion matrix

In []: yhat = tree_cv.predict(X_test)

did not land

TASK 10

tree = DecisionTreeClassifier()

tree_cv.fit(X_train, Y_train)

accuracy: 0.8875

In []: Y_hat = tree_cv.predict(Y_test)

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

accuracy_score()

TASK 9

did not land

True labels

Out[]: •

'max_features': ['auto', 'sqrt'], 'min_samples_leaf': [1, 2, 4], 'min_samples_split': [2, 5, 10]}

'max_depth': [2*n for n in range(1,10)],

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

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

- 12

'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute'],

In []: print("tuned hpyerparameters :(best parameters) ",knn_cv.best_params_) #tune出best parameter

- 12

- 10

Date (YYYY-MM-DD) Version

1.0

2022-11-09

Changed By

IBM Corporation 2022. All rights reserved.

Change Description

Pratiksha Verma Converted initial version to Jupyterlite

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

12

land

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

12

plt.figure(figsize=(10, 6))

plt.xlabel('Metrics', fontsize=14) plt.ylabel('Score', fontsize=14)

- 12

- 10

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

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.

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

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.

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.

tuned hpyerparameters :(best parameters) {'criterion': 'entropy', 'max_depth': 6, 'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_split': 10, 'splitter': 'best'}

In []: log_precision = precision_score(Y_test, Y_hat) log_recall = recall_score(Y_test, Y_hat)

0.8333333333333334 0.8 1.0 0.8888888888888888

Confusion Matrix

Predicted labels

log_f1 = f1_score(Y_test, Y_hat)

plot_confusion_matrix(Y_test, yhat)

Lets look at the confusion matrix:

In []: yhat=logreg_cv.predict(X_test)

did not land

TASK 6

svm = SVC()

In []: svm_cv.fit(X_train, Y_train)

▶ estimator: SVC

▶ SVC

In []: svm_cv.best_estimator_

TASK 7

Out[]: ▼

tuned hpyerparameters :(best parameters) {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}

text2 = io.BytesIO((await resp2.arrayBuffer()).to_py())

df.to_csv('dataset_part_2.csv', index= None)

1 2010-06-04

2 2012-05-22

3 2013-03-01

4 2013-09-29

5 2013-12-03

6104.959412

525.000000

677.000000

500.000000

3170.000000

86.0 15400.000000

87.0 15400.000000

88.0 15400.000000

89.0 15400.000000

3681.000000

from sklearn.preprocessing import StandardScaler

hyperparameters are selected using the function GridSearchCV.

X = StandardScaler().fit_transform(X)

X_train, X_test, Y_train, Y_test

we can see we only have 18 test samples.

resp2 = await fetch(URL2)

 $X = pd.read_csv(text2)$

 $X = pd.read_csv(URL2)$

FlightNumber

2.0

3.0

4.0

90.0

90 rows × 83 columns

TASK 1

Y = df['Class']

TASK 2

In []: # students get this

TASK 3

In []: len(df)

In []: len(Y_test)

len(X_test)

TASK 4

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

best_score_ .

TASK 5

In []: # Calculate accuracy on test data

logreg=LogisticRegression()

logreg_cv.fit(X_train, Y_train)

GridSearchCV

▶ estimator: LogisticRegression

▶ LogisticRegression

print("accuracy :", logreg_cv.best_score_)

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

log_accuracy = logreg_cv.score(X_test, Y_test) print("Accuracy on test data:", log_accuracy)

log_accuracy = accuracy_score(Y_test, Y_hat) log_precision = precision_score(Y_test, Y_hat)

data = {'metrics': metrics, 'score': score}

log_recall = recall_score(Y_test, Y_hat)

Accuracy on test data: 0.8333333333333333

Y_hat = logreg_cv.predict(X_test)

log_f1 = f1_score(Y_test, Y_hat)

df_sorted = pd.DataFrame(data)

plt.title('classification metrics')

plt.xlabel('Metrics') plt.ylabel('Score')

plt.xticks(rotation=45)

plt.show()

1.0

0.8

0.4

0.2

did not land

True labels

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

from sklearn.model_selection import GridSearchCV

logreg_cv =GridSearchCV(logreg, parameters, cv=10)

Out[]: 90

Out[]: 18

Out[]: ▶

Y = np.asanyarray(Y)

X = np.asanyarray(X)

cm = confusion_matrix(y, y_predict)

ax.set_xlabel('Predicted labels')

ax.set_ylabel('True labels') ax.set_title('Confusion Matrix');

Decision Tree classification algorithm

This function is to plot the confusion matrix.

In []: def plot_confusion_matrix(y,y_predict):

ax= plt.subplot()

Load the dataframe

resp1 = await fetch(URL1)

data = pd.read_csv(text1)

df = pd.read_csv(URL1)

plt.show()

Load the data

In []: from js import fetch import io

df.head()

0

2

4

In []: X

0

2

3

4

85

86

87

88

89

Out[]:

FlightNumber

In []:

Out[]:

from sklearn.model_selection import train_test_split

Allows us to split our data into training and testing data

Allows us to test parameters of classification algorithms and find the best one

sns.heatmap(cm, annot=True, ax = ax); #annot=True to annotate cells

ax.xaxis.set_ticklabels(['did not land', 'land']); ax.yaxis.set_ticklabels(['did not land', 'landed'])

URL1 = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_2.csv"

In []: URL1 = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_2.csv"

LaunchSite

None None

None None

None None

None None

VAFB SLC 4E False Ocean

URL2 = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_3.csv'

URL2 = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321EN-SkillsNetwork/datasets/dataset_part_3.csv'

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

0.0

1.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

0.0

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']).

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

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

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.

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

6104.959412 LEO CCAFS SLC 40

525.000000 LEO CCAFS SLC 40

PO

Falcon 9 3170.000000 GTO CCAFS SLC 40

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

0.0

ISS CCAFS SLC 40

Outcome Flights GridFins Reused Legs LandingPad Block ReusedCount Serial

NaN

NaN

NaN

NaN

NaN

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 Reu

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

0.0

0.0

0.0

1.0

0.0

0.0

0.0

False 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 ...

False

False

False

False

False

Longitude

-80.577366 28.561857

-80.577366 28.561857

-80.577366 28.561857

-120.610829 34.632093

-80.577366 28.561857

0 B0003

0 B0005

0 B0007

0 B1003

0 B1004

0.0

0.0

0.0

0.0

0.0

1.0

0.0

0.0

1.0

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

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

Latitude Class

0

0

0

0

0

• Find the method performs best using test data

Standardize the data

Objectives

In []: import piplite

lab, you will create a machine learning pipeline to predict if the first stage will land given the data from the preceding labs.

Estimated time needed: 60 minutes



Skills Space X Falcon 9 First Stage Landing Prediction

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

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 mathematic

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

Network