



- ✓ **Space X Falcon 9 First Stage Landing Prediction**
- ✓ Assignment: Machine Learning Prediction

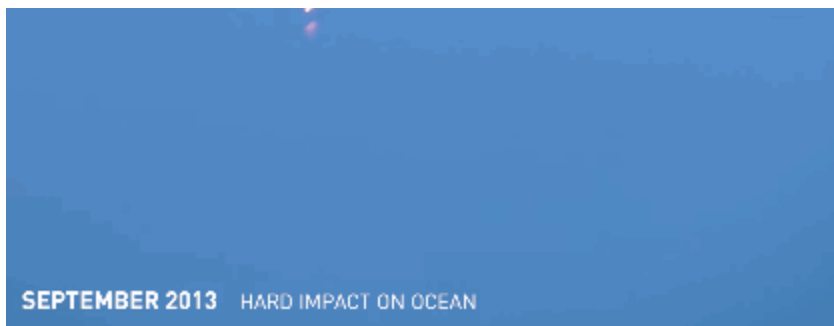
Estimated time needed: **60** minutes

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 Libraries and Define Auxiliary Functions

```
1 import piplite
2 await piplite.install(['numpy'])
3 await piplite.install(['pandas'])
4 await piplite.install(['seaborn'])
```

We will import the following libraries for the lab

```
1 # Pandas is a software library written for the Python programming language for data mani
2 import pandas as pd
3 # NumPy is a library for the Python programming language, adding support for large, multi
4 import numpy as np
5 # Matplotlib is a plotting library for python and pyplot gives us a MatLab like plotting
6 import matplotlib.pyplot as plt
7 #Seaborn is a Python data visualization library based on matplotlib. It provides a high-
8 import seaborn as sns
9 # Preprocessing allows us to standarsize our data
10 from sklearn import preprocessing
```

```

10 from sklearn import preprocessing
11 # Allows us to split our data into training and testing data
12 from sklearn.model_selection import train_test_split
13 # Allows us to test parameters of classification algorithms and find the best one
14 from sklearn.model_selection import GridSearchCV
15 # Logistic Regression classification algorithm
16 from sklearn.linear_model import LogisticRegression
17 # Support Vector Machine classification algorithm
18 from sklearn.svm import SVC
19 # Decision Tree classification algorithm
20 from sklearn.tree import DecisionTreeClassifier
21 # K Nearest Neighbors classification algorithm
22 from sklearn.neighbors import KNeighborsClassifier

```

This function is to plot the confusion matrix.

```

1 def plot_confusion_matrix(y,y_predict):
2     "this function plots the confusion matrix"
3     from sklearn.metrics import confusion_matrix
4
5     cm = confusion_matrix(y, y_predict)
6     ax= plt.subplot()
7     sns.heatmap(cm, annot=True, ax = ax); #annot=True to annotate cells
8     ax.set_xlabel('Predicted labels')
9     ax.set_ylabel('True labels')
10    ax.set_title('Confusion Matrix');
11    ax.xaxis.set_ticklabels(['did not land', 'land']); ax.yaxis.set_ticklabels(['did n
12    plt.show()

```

✓ Load the dataframe

Load the data

```

1 from js import fetch
2 import io
3
4 URL1 = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321E
5 resp1 = await fetch(URL1)
6 text1 = io.BytesIO((await resp1.arrayBuffer()).to_py())
7 data = pd.read_csv(text1)

```



```

1 data.head()

```

	FlightNumber	Date	BoosterVersion	PayloadMass	Orbit	LaunchSite	Outcome
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

```

1 URL2 = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DS0321E
2 resp2 = await fetch(URL2)
3 text2 = io.BytesIO((await resp2.arrayBuffer())).to_py()
4 X = pd.read_csv(text2)

```

```
1 X.head(100)
```

	FlightNumber	PayloadMass	Flights	Block	ReusedCount	Orbit_ES- L1	Orbit_GEO	Ort
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']`).

```
1 import pandas as pd
2
3 Y = data['Class'].to_numpy()
```

✓ TASK 2

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

```
1 # students get this
2 transform = preprocessing.StandardScaler()
3 X = preprocessing.StandardScaler().fit(X).transform(X.astype(float))
```

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`

```
1 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.

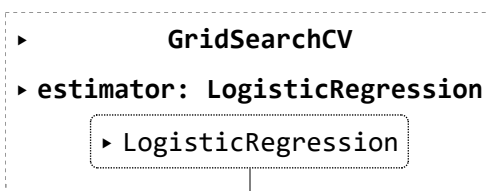
```
1 Y_test.shape
(18,)
```

✓ 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`.

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

1 parameters = {"C":[0.01,0.1,1], 'penalty':['l2'], 'solver':['lbfgs']}# l1 lasso l2 ridge
2 lr=LogisticRegression()
3
4 # Create a GridSearchCV object with cv = 10
5 logreg_cv = GridSearchCV(lr, parameters, cv=10)
6
7 # Fit the GridSearchCV object to the training data
8 logreg_cv.fit(X_train, Y_train)
9
```



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_`.

```
1 print("tuned hpyerparameters :(best parameters) ",logreg_cv.best_params_)
2 print("accuracy :",logreg_cv.best_score_)

tuned hpyerparameters :(best parameters) {'C': 0.01, 'penalty': 'l2', 'solver': 'lbf
accuracy : 0.8464285714285713
```

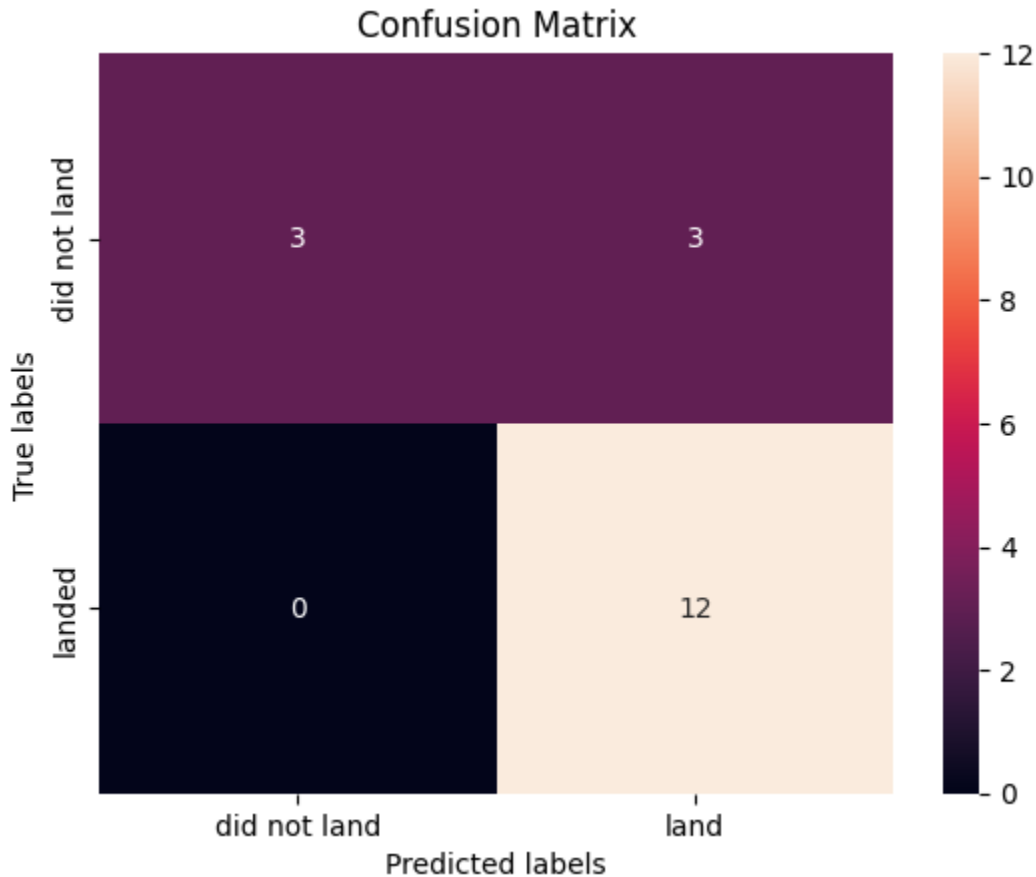
✓ TASK 5

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

```
1 test_accuracy = logreg_cv.score(X_test, Y_test)
2
```

Lets look at the confusion matrix:

```
1 yhat=logreg_cv.predict(X_test)
2 plot_confusion_matrix(Y_test,yhat)
```



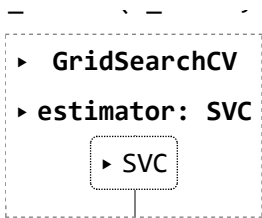
Examining the confusion matrix, we see that logistic regression can distinguish between the different classes. We see that the major problem is false positives.

✓ 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`.

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

1 # Create a GridSearchCV object with cv = 10
2 svm_cv = GridSearchCV(svm, parameters, cv=10)
3
4 # Fit the GridSearchCV object to the training data
5 svm_cv.fit(X_train, Y_train)
```



```

1 print("tuned hpyerparameters :(best parameters) ",svm_cv.best_params_)
2 print("accuracy :",svm_cv.best_score_)

tuned hpyerparameters :(best parameters) {'C': 1.0, 'gamma': 0.03162277660168379, 'k
accuracy : 0.8482142857142856

```

✓ TASK 7

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

```

1 test_accuracy = svm_cv.score(X_test, Y_test)
2

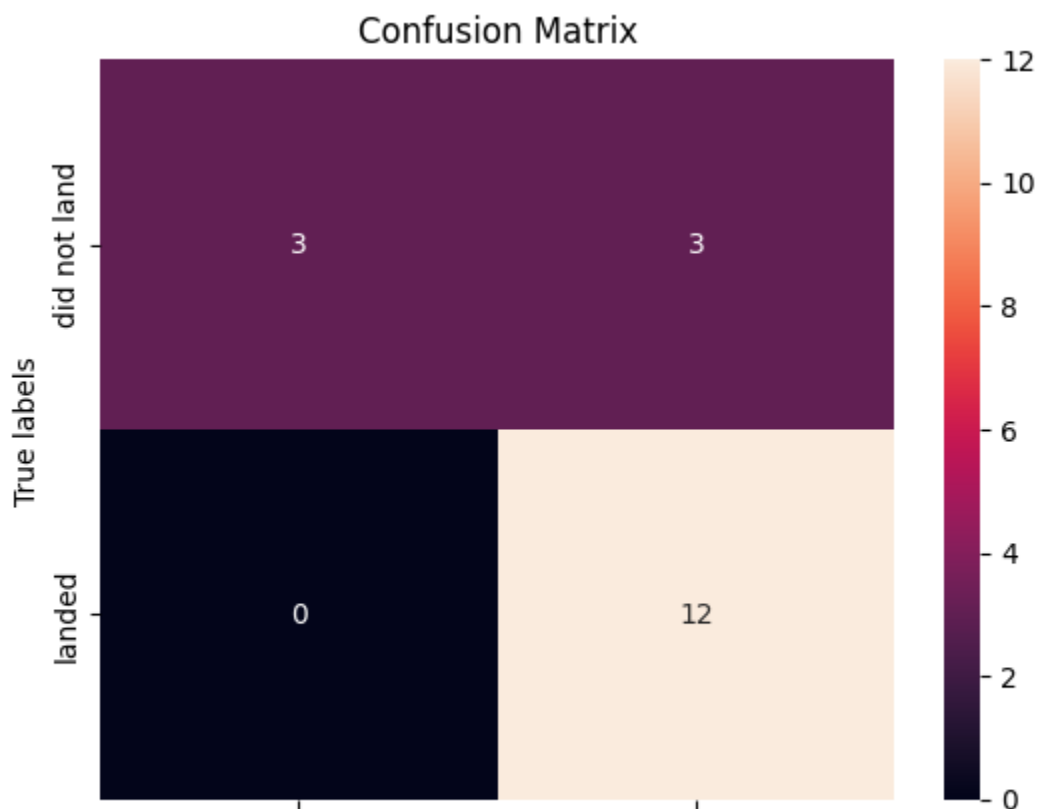
```

We can plot the confusion matrix

```

1 yhat=svm_cv.predict(X_test)
2 plot_confusion_matrix(Y_test,yhat)

```



did not land land

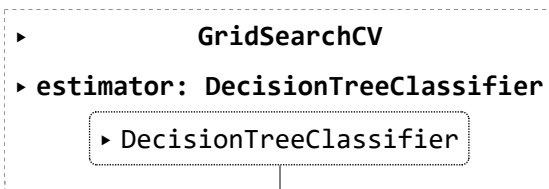
Predicted labels

✓ 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`.

```
1 parameters = {
2     'criterion': ['gini', 'entropy'],
3     'splitter': ['best', 'random'],
4     'max_depth': [2*n for n in range(1,10)],
5     'max_features': ['auto', 'sqrt'],
6     'min_samples_leaf': [1, 2, 4],
7     'min_samples_split': [2, 5, 10]
8 }
9
10 tree = DecisionTreeClassifier()

1 # Create a GridSearchCV object with cv = 10
2 tree_cv = GridSearchCV(tree, parameters, cv=10)
3
4 # Fit the GridSearchCV object to the training data
5 tree_cv.fit(X_train, Y_train)
```



```
1 print("tuned hpyerparameters :(best parameters) ",tree_cv.best_params_)
2 print("accuracy :",tree_cv.best_score_)

tuned hpyerparameters :(best parameters) {'criterion': 'gini', 'max_depth': 6, 'max_
accuracy : 0.8767857142857143
```

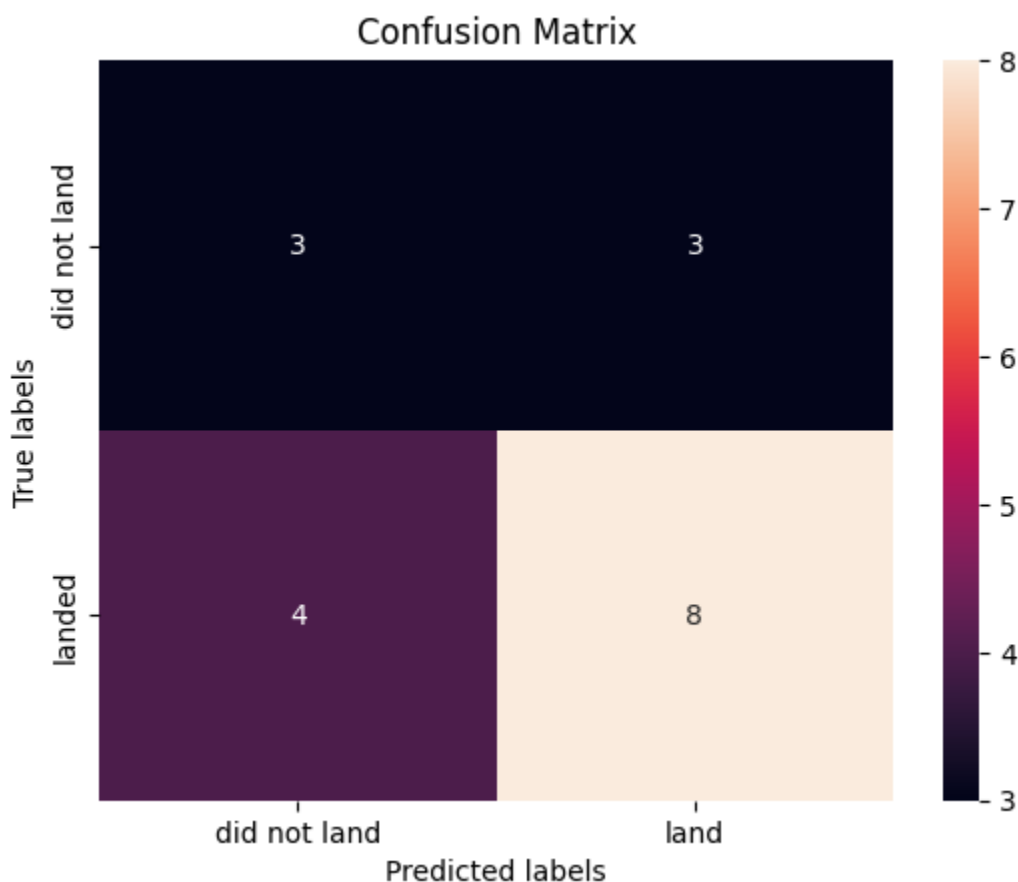
✓ TASK 9

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

```
1 test_accuracy = tree_cv.score(X_test, Y_test)
```

We can plot the confusion matrix

```
1 yhat = tree_cv.predict(X_test)
2 plot_confusion_matrix(Y_test,yhat)
```



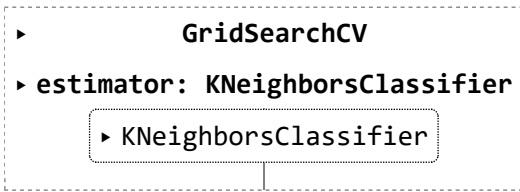
✓ TASK 10

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`.

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

1 # Create a GridSearchCV object with cv = 10
2 KNN_cv = GridSearchCV(KNN, parameters, cv=10)
3
4 # Fit the GridSearchCV object to the training data
5 KNN_cv.fit(X_train, y_train)
```

```
5 KNN_cv.fit(X_train, Y_train)
6
```



```
1 print("tuned hpyerparameters :(best parameters) ",KNN_cv.best_params_)
2 print("accuracy :",KNN_cv.best_score_)

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

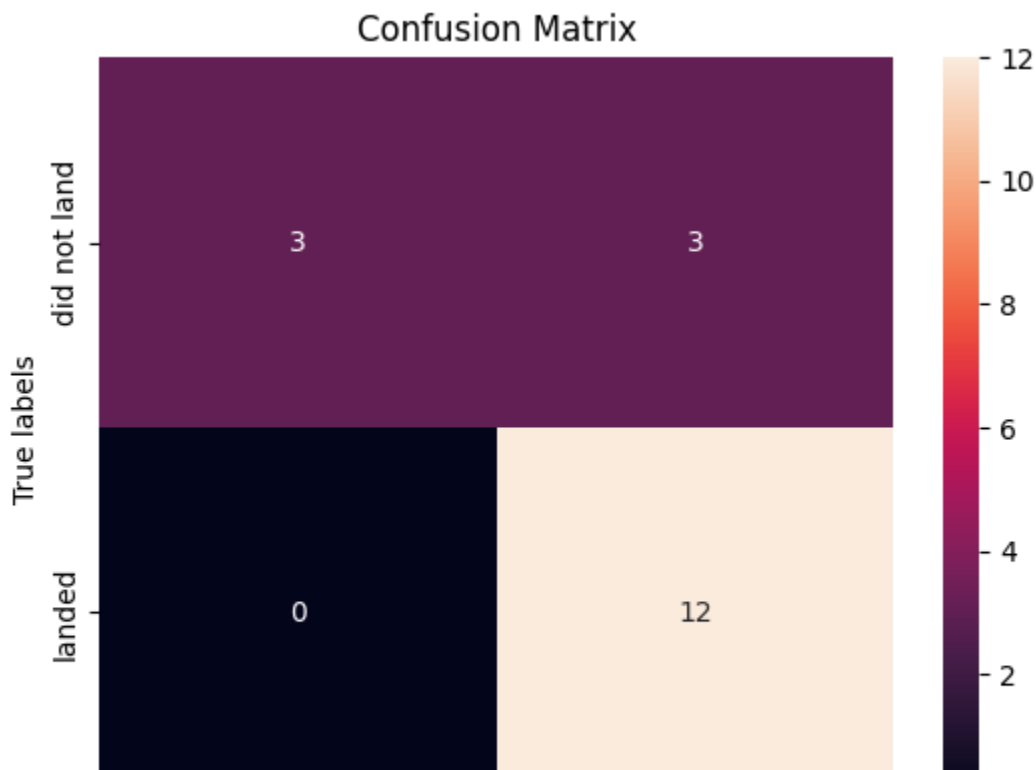
✓ TASK 11

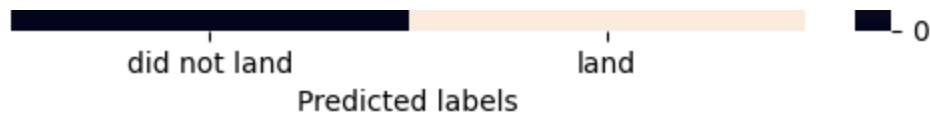
Calculate the accuracy of knn_cv on the test data using the method `score` :

```
1 test_accuracy = KNN_cv.score(X_test, Y_test)
```

We can plot the confusion matrix

```
1 yhat = KNN_cv.predict(X_test)
2 plot_confusion_matrix(Y_test,yhat)
```





✓ TASK 12

Find the method performs best:

```
1 # Compare test accuracies
2 model_accuracies = {
3     'Logistic Regression': test_accuracy,
4     'Support Vector Machine': svm_test_accuracy,
5     'Decision Tree': tree_test_accuracy,
6     'k-Nearest Neighbors': knn_test_accuracy
7 }
8
9 best_model = max(model_accuracies, key=model_accuracies.get)
10 print("Best Model:", best_model)
11 print("Best Model Test Accuracy:", model_accuracies[best_model])
```

NameError Traceback (most recent call last)

Cell In[47], line 4

```
1 # Compare test accuracies
2 model_accuracies = {
3     'Logistic Regression': test_accuracy,
----> 4     'Support Vector Machine': svm_test_accuracy,
5     'Decision Tree': tree_test_accuracy,
6     'k-Nearest Neighbors': knn_test_accuracy
7 }
8
9 best_model = max(model_accuracies, key=model_accuracies.get)
10 print("Best Model:", best_model)
```

NameError: name 'svm_test_accuracy' is not defined

✓ 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

IBM Corporation 2022. All rights reserved.