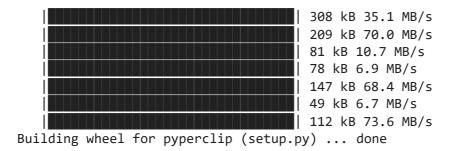
```
from google.colab import drive
drive.mount('/content/gdrive')
          Mounted at /content/gdrive

root_path = "gdrive/MyDrive/ML/"
```

!pip install --quiet optuna



```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import matplotlib
import seaborn as sns
from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
from sklearn.preprocessing import OrdinalEncoder, OneHotEncoder, StandardScaler
from sklearn.linear_model import LinearRegression
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import KNNImputer, IterativeImputer
from sklearn.naive_bayes import MultinomialNB
from sklearn.model_selection import cross_val_score, train_test_split
from sklearn.feature_selection import mutual_info_classif
from sklearn.mixture import GaussianMixture
from sklearn.cluster import KMeans, DBSCAN
from yellowbrick.cluster import KElbowVisualizer
from sklearn.svm import SVC
import optuna
```

Data Loading

```
df = pd.read_csv(root_path + "train.csv", index_col='PassengerId').reset_index(drop=True)
df_test = pd.read_csv(root_path + "test.csv", index_col='PassengerId').reset_index(drop=Tr
df.head()
```

9/22, 10:17	spacesnip-titanic-comp-u-80219-acc.lpynb - Colaboratory										
	HomePlanet	CryoSleep	Cabin	Destination	Age	VIP	RoomService	FoodCourt	Sho		
	0 Europa	ı False	B/0/P	TRAPPIST- 1e	39.0	False	0.0	0.0			
	1 Earth	ı False	F/0/S	TRAPPIST- 1e	24.0	False	109.0	9.0			
print(print("Shape ")) "Train: ", df "Test: ", df_										
9	Shape										
	rain: (8693, est: (4277,	•									
print(print(print("NaN values ') "Train: \n", "Test: \n", o	df.isna().su		·							
C C C C C F F F S S V N	Train: HomePlanet CryoSleep Cabin Destination Age VIP RoomService FoodCourt ShoppingMall Spa VRDeck Lame Transported Htype: int64	201 217 199 182 179 203 181 183 208 183 188 200									
C C E A	Test: HomePlanet CryoSleep Cabin Destination Age VIP RoomService	87 93 100 92 91 93 82									

df["Transported"] = df["Transported"].astvpe(np.int8) $https://colab.research.google.com/drive/1Am-BsEJ9RWTTb7fxM-4sn_dCn3c_YVeE\#scrollTo=5l-r0ohciNc0\&printMode=true$

FoodCourt

Spa

Name

VRDeck

ShoppingMall

dtype: int64

106

98

101

80 94 df.head()

	HomePlanet	CryoSleep	Cabin	Destination	Age	VIP	RoomService	FoodCourt	Sho
0	Europa	False	B/0/P	TRAPPIST- 1e	39.0	False	0.0	0.0	
1	Earth	False	F/0/S	TRAPPIST- 1e	24.0	False	109.0	9.0	
2	Europa	False	A/0/S	TRAPPIST- 1e	58.0	True	43.0	3576.0	
4									•

df.Transported.value_counts()

1 4378

0 4315

Name: Transported, dtype: int64

cabin_splited = df["Cabin"].str.split("/", expand=True)
cabin_splited.columns = ["Cabin_deck", "Cabin_num", "Cabin_side"]
cabin_splited.head()

	Cabin_deck	Cabin_num	Cabin_side	1
0	В	0	Р	
1	F	0	S	
2	А	0	S	
3	А	0	S	
4	F	1	S	

df1 = pd.concat([df, cabin_splited], axis=1).drop("Cabin", axis=1)
df1.head()

	HomePlanet	CryoSleep	Destination	Age	VIP	RoomService	FoodCourt	ShoppingMa
0	Europa	False	TRAPPIST- 1e	39.0	False	0.0	0.0	
1	Earth	False	TRAPPIST- 1e	24.0	False	109.0	9.0	2
2	Europa	False	TRAPPIST- 1e	58.0	True	43.0	3576.0	
4								>

cabin_splited = df_test["Cabin"].str.split("/", expand=True)
cabin_splited.columns = ["Cabin_deck", "Cabin_num", "Cabin_side"]
df_test_1 = pd.concat([df_test, cabin_splited], axis=1).drop("Cabin", axis=1)
df_test_1.head()

	HomePlanet	CryoSleep	Destination	Age	VIP	RoomService	FoodCourt	ShoppingMa
0	Earth	True	TRAPPIST- 1e	27.0	False	0.0	0.0	
1	Earth	False	TRAPPIST- 1e	19.0	False	0.0	9.0	
2	Europa	True	55 Cancri e	31.0	False	0.0	0.0	
4								>

```
object_cols = [i for i in df1.columns if df1[i].dtype == "0"]
for i in object_cols:
    print(i, ": ", df1[i].nunique())
    HomePlanet : 3
```

Destination : 3
VIP : 2
Name : 8473
Cabin_deck : 8
Cabin_num : 1817
Cabin_side : 2

CryoSleep: 2

Since Cabin_num hasn't any relation ship is a passenger Transported, we'll remove it. With Name column the same thing

```
df2 = df1.drop(["Name", "Cabin_num"], axis=1)
df_test_2 = df_test_1.drop(["Name", "Cabin_num"], axis=1)
df2.head()
```

	HomePlanet	CryoSleep	Destination	Age	VIP	RoomService	FoodCourt	ShoppingMa
0	Europa	False	TRAPPIST- 1e	39.0	False	0.0	0.0	
1	Earth	False	TRAPPIST- 1e	24.0	False	109.0	9.0	2
2	Europa	False	TRAPPIST- 1e	58.0	True	43.0	3576.0	
4								>

Handle NaN values

[] Ļ9 cells hidden

Encoding values

Model Creation and Training

```
model = RandomForestClassifier(n estimators=100, random state=0)
scores = cross_val_score(model, X, y, cv=10)
scores
     array([0.79770115, 0.75402299, 0.78390805, 0.79171461, 0.80437284,
            0.81127733, 0.80552359, 0.79056387, 0.80782509, 0.79516686])
def run(trial):
    criterion = trial.suggest_categorical('criterion', ['gini', 'entropy'])
    bootstrap = trial.suggest_categorical('bootstrap',['True','False'])
    max_depth = trial.suggest_int('max_depth', 2, 50)
    max_features = trial.suggest_categorical('max_features', ['auto', 'sqrt','log2'])
    max_leaf_nodes = trial.suggest_int('max_leaf_nodes', 2, 100)
    n_estimators = trial.suggest_int('n_estimators', 100, 500)
    X_train, X_val, y_train, y_val = train_test_split(X, y)
    regr = RandomForestClassifier(bootstrap = bootstrap, criterion = criterion,
                                  max_depth = max_depth, max_features = max_features,
                                  max_leaf_nodes = max_leaf_nodes,n_estimators = n_estimato
    regr.fit(X_train, y_train)
    return regr.score(X_val, y_val)
study = optuna.create_study(direction='maximize')
study.optimize(run, n_trials=200)
print(study.best_params)
     poolstrap : raise , max_depth : 41, max_teatures : logz , max_leat_hodes : //
    'bootstrap': 'False', 'max_depth': 43, 'max_features': 'log2', 'max_leaf_nodes': 77
    'bootstrap': 'False', 'max_depth': 41, 'max_features': 'log2', 'max_leaf_nodes': 74
    'bootstrap': 'False', 'max_depth': 25, 'max_features': 'log2', 'max_leaf_nodes': 81
    'bootstrap': 'False', 'max_depth': 6, 'max_features': 'log2', 'max_leaf_nodes': 84,
    bootstrap': 'False', 'max_depth': 45, 'max_features': 'log2', 'max_leaf_nodes': 79, bootstrap': 'False', 'max_depth': 40, 'max_features': 'log2', 'max_leaf_nodes': 92,
    , 'bootstrap': 'False', 'max_depth': 37, 'max_features': 'sqrt', 'max_leaf_nodes':
    'bootstrap': 'True', 'max_depth': 9, 'max_features': 'auto', 'max_leaf_nodes': 100,
    ', 'bootstrap': 'False', 'max_depth': 39, 'max_features': 'auto', 'max_leaf_nodes':
    'bootstrap': 'False', 'max_depth': 35, 'max_features': 'sqrt', 'max_leaf_nodes': 75
    'bootstrap': 'True', 'max_depth': 8, 'max_features': 'auto', 'max_leaf_nodes': 97,
    'bootstrap': 'False', 'max_depth': 40, 'max_features': 'log2', 'max_leaf_nodes': 87
    'bootstrap': 'True', 'max_depth': 12, 'max_features': 'auto', 'max_leaf_nodes': 94,
    'bootstrap': 'True', 'max_depth': 15, 'max_features': 'auto', 'max_leaf_nodes': 91,
    bootstrap': 'True', 'max_depth': 10, 'max_features': 'auto', 'max_leaf_nodes': 7,
    , 'bootstrap': 'False', 'max_depth': 38, 'max_features': 'sqrt', 'max_leaf_nodes':
    'bootstrap': 'True', 'max_depth': 30, 'max_features': 'auto', 'max_leaf_nodes': 95,
    'bootstrap': 'True', 'max_depth': 42, 'max_features': 'auto', 'max_leaf_nodes': 98,
    bootstrap': 'False', 'max_depth': 7, 'max_features': 'auto', 'max_leaf_nodes': 61,
```

```
spaceship-titanic-comp-0-80219-acc.ipynb - Colaboratory
'bootstrap': 'True', 'max_depth': 11, 'max_features': 'auto', 'max_leaf_nodes': 98,
bootstrap': 'False', 'max depth': 36, 'max features': 'auto', 'max leaf nodes': 65,
, 'bootstrap': 'True', 'max_depth': 20, 'max_features': 'auto', 'max_leaf_nodes': 9
'bootstrap': 'True', 'max_depth': 13, 'max_features': 'sqrt', 'max_leaf_nodes': 76,
bootstrap': 'False', 'max depth': 32, 'max features': 'auto', 'max leaf nodes': 106
'bootstrap': 'False', 'max_depth': 30, 'max_features': 'auto', 'max_leaf_nodes': 96
', 'bootstrap': 'False', 'max_depth': 32, 'max_features': 'log2', 'max_leaf_nodes':
, 'bootstrap': 'False', 'max_depth': 32, 'max_features': 'log2', 'max_leaf nodes':
'bootstrap': 'False', 'max_depth': 32, 'max_features': 'log2', 'max_leaf_nodes': 10
'bootstrap': 'False', 'max_depth': 34, 'max_features': 'log2', 'max_leaf_nodes': 16
'bootstrap': 'False', 'max_depth': 29, 'max_features': 'auto', 'max_leaf_nodes': 98
'bootstrap': 'False', 'max_depth': 29, 'max_features': 'log2', 'max_leaf_nodes': 10
'bootstrap': 'False', 'max_depth': 28, 'max_features': 'auto', 'max_leaf_nodes': 94
'bootstrap': 'False', 'max_depth': 31, 'max_features': 'auto', 'max_leaf_nodes': 97
'bootstrap': 'False', 'max_depth': 27, 'max_features': 'auto', 'max_leaf_nodes': 91
bootstrap': 'False', 'max_depth': 26, 'max_features': 'log2', 'max_leaf_nodes': 95,
'bootstrap': 'False', 'max_depth': 31, 'max_features': 'auto', 'max_leaf_nodes': 9
'bootstrap': 'False', 'max_depth': 30, 'max_features': 'auto', 'max_leaf_nodes': 93
'bootstrap': 'False', 'max_depth': 33, 'max_features': 'auto', 'max_leaf_nodes': 89
'bootstrap': 'False', 'max_depth': 32, 'max_features': 'auto', 'max_leaf_nodes': 89' bootstrap': 'False', 'max_depth': 33, 'max_features': 'auto', 'max_leaf_nodes': 98'
'bootstrap': 'False', 'max_depth': 31, 'max_features': 'auto', 'max_leaf_nodes': 92 'bootstrap': 'False', 'max_depth': 29, 'max_features': 'auto', 'max_leaf_nodes': 88
bootstrap': 'False', 'max_depth': 34, 'max_features': 'auto', 'max_leaf_nodes': 96,
', 'bootstrap': 'False', 'max_depth': 33, 'max_features': 'log2', 'max_leaf_nodes':
'bootstrap': 'False', 'max_depth': 28, 'max_features': 'auto', 'max_leaf_nodes': 94
'bootstrap': 'True', 'max_depth': 32, 'max_features': 'auto', 'max_leaf_nodes': 98,
'bootstrap': 'False', 'max_depth': 30, 'max_features': 'auto', 'max_leaf_nodes': 10
'bootstrap': 'False', 'max_depth': 31, 'max_features': 'auto', 'max_leaf_nodes': 86
bootstrap': 'False', 'max depth': 33, 'max features': 'auto', 'max leaf nodes': 91,
'bootstrap': 'False', 'max_depth': 30, 'max_features': 'auto', 'max_leaf_nodes': 96
'bootstrap': 'False', 'max_depth': 9, 'max_features': 'auto', 'max_leaf_nodes': 81,
'bootstrap': 'False', 'max_depth': 34, 'max_features': 'log2', 'max_leaf_nodes': 53
'bootstrap': 'True', 'max_depth': 27, 'max_features': 'auto', 'max_leaf_nodes': 78,
', 'bootstrap': 'False', 'max_depth': 35, 'max_features': 'auto', 'max_leaf_nodes':
'bootstrap': 'True', 'max_depth': 29, 'max_features': 'auto', 'max_leaf_nodes': 98,
tors': 465}
 {'criterion': 'gini',
```

```
rf_params = study.best_params
rf params
      'bootstrap': 'True',
      'max_depth': 37,
      'max_features': 'auto',
      'max leaf nodes': 100,
      'n estimators': 465}
model = RandomForestClassifier(**rf params, n jobs=2, random state=0)
model.fit(X, y)
print("Done!")
     Done!
```

	PassengerId	Transported	
0	0013_01	True	
1	0018_01	False	
2	0019_01	True	
3	0021_01	True	
4	0023_01	True	
4272	9266_02	True	
4273	9269_01	False	
4274	9271_01	True	
4275	9273_01	True	
4276	9277_01	True	

ss.to_csv(root_path + "submission.csv", index=False)

4277 rows × 2 columns

model.fit(X_train, y_train)

```
X_train, X_val, y_train, y_val = train_test_split(X, y)
from scipy.sparse.sputils import matrix
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix
model = LogisticRegression(max iter = 1000)
```

prediction = model.predict(X_val)
print('The accuracy of the Logistic Regression is', accuracy_score(prediction, y_val))

The accuracy of the Logistic Regression is 0.7805887764489421 /usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py:818: Conver STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,

completed at 8:47 PM

✓ 2s

×