

UCK 358E: Project #1 21600 - Monday

Due on April 30, 2023

Asst. Assoc. Prof. Dr. Barış Başpınar

Hüseyin Tutan * - 110190021

^{*}tutan 19@itu.edu.tr

Contents

1	Analysis of the Given Code	1
	1.1 Control of Overfitting or Underfitting	2
2	Data Preprocessing	5
	2.1 Missing Values	5
	2.2 Feature Engineering	6
3	B Different Machine Learning Models	6
	3.1 GradientBoostingClassifier	6
	3.2 Support Vector Machine	7
	3.3 Naive Bayes	7
	3.4 Logistic Regression	7
4	Hyperparameter Tuning for the Best Scored ML Method	8
5	Final Python Code	9

1 Analysis of the Given Code

First, let's make the code executable and add some basic comment lines.

```
#!/usr/bin/env python3
  # @Author: Huseyin Tutan
  # Import the libraries
  import pandas as pd
  import numpy as np
  import matplotlib.pyplot as plt
11
  # Import the data
13
14 data = pd.read_csv("train.csv")
  test = pd.read_csv("test.csv")
16
  # Code Section
18
  y = data["Survived"]
  features = ["Pclass", "Sex", "SibSp", "Parch"]
  X = pd.get_dummies(data[features])
  X_test = pd.get_dummies(test[features])
23
  model = RandomForestClassifier(n_estimators=100, max_depth=5, random_state=1)
26 model.fit(X, y)
27 predictions = model.predict(X_test)
28 output = pd.DataFrame({'PassengerId': test.PassengerId, 'Survived': predictions})
29 output.to_csv('submission1.csv', index=False)
```

1.1 Control of Overfitting or Underfitting

Firstly, I planned to check is there any suitability problems. Then, I did a quick search on kaggle.com discussion forum and found that post. "Overfitting and underfitting the Titanic" I adapted the code to my own code for check our code's fitting status and our code last section became like this.

```
1 # read in the data
2 train_data = pd.read_csv('train.csv')
3 test_data = pd.read_csv('test.csv')
            = pd.read_csv('submission1.csv')
  solution
  # select some features
  features = ["Pclass", "Sex", "SibSp", "Parch"]
          = pd.get_dummies(train_data[features])
          = train_data["Survived"]
  final_X_test = pd.get_dummies(test_data[features])
11
  # perform the classification and the fit
  classifier = RandomForestClassifier(random_state=4)
  classifier.fit(X, y)
15
predictions = classifier.predict(final_X_test)
  K_{splits} = 11
18
  # calculate the scores
  CV_scores = cross_val_score(classifier, X, y, cv=K_splits)
  # Print Section & Visualize
22
23
  print ("The mean accuracy score of the train data is %.5f" % classifier.score (X, y))
  print ("The mean cross-validation score is %.5f \pm %0.2f" % (CV_scores.mean (), CV_score
  print("The test (i.e. leaderboard) score is %.5f" % accuracy_score(solution[|'Survived'],
27
 print ("The individual cross-validation scores are: \n", CV_scores)
29 print("The minimum cross-validation score is %.3f" % min(CV_scores))
  print("The maximum cross-validation score is %.3f" % max(CV_scores))
31
32 cv = StratifiedKFold(n_splits=K_splits)
33 visualizer = CVScores(classifier, cv=cv, scoring='f1_weighted', size=(1200, 400))
34 visualizer.fit(X, y)
35 visualizer.show()
```

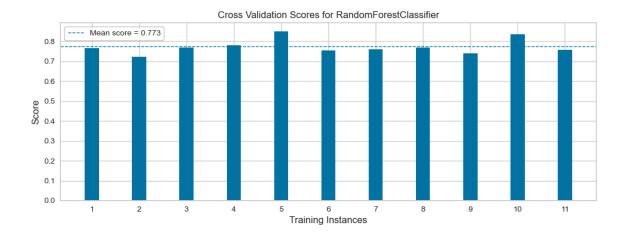
Also added these headers to code:

```
import colorama

from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from colorama import Fore
from sklearn.ensemble import RandomForestClassifier
from yellowbrick.model_selection import CVScores
from sklearn.model_selection import StratifiedKFold
```

Finally, the terminal output is as follows:

```
1 The mean accuracy score of the train data is 0.81706
2 The mean cross-validation score is 0.77890 \pm 0.07
3 The test (i.e. leaderboard) score is 0.95215
4 The individual cross-validation scores are:
5 [0.7654321 0.72839506 0.77777778 0.7777778 0.85185185 0.7654321
6 0.7654321 0.77777778 0.75308642 0.83950617 0.7654321 ]
7 The minimum cross-validation score is 0.728
8 The maximum cross-validation score is 0.852
```



Direct Quote from web-site:

"We now find that the

That said, it is more usual to find that

$$CV(score) \approx LB(score)$$

up to the point of overfitting. Also note that the leaderboard score is lower when underfitting (0.66) than it was in the overfitting case above (0.76)."

So, we can say there is no over or underfitting issues in our code.

2 Data Preprocessing

In this section, after analyzing the code, I reported the different methods I tried in the preprocessing phase.

2.1 Missing Values

Firstly, I analyzed is there any missing values in our data with this code.

```
print(data.isna().sum())
```

Terminal output looks like below:

```
0
1 PassengerId
2 Survived
                      0
3 Pclass
                     0
4 Name
  Sex
                     0
6 Age
                   177
                     0
  SibSp
  Parch
                     0
 Ticket
                     0
10 Fare
                     0
11 Cabin
                    687
12 Embarked
                     2
13 dtype: int64
```

As seen in terminal, we have some missing values at Ages, Cabin and Embarked data. Afterwards, I completed the missing data with the following codes and the specified methods. I filled the age part with the mean method

```
1 data['Age'].fillna(data['Age'].mean(), inplace=True)
```

I filled the embarked part with the mod method

```
1 data['Embarked'].fillna(data['Embarked'].mode()[0], inplace=True)
```

In cabin column there are so many missing values so I decided to do not complete it and do not use it for training my model.

2.2 Feature Engineering

I did research about some useful patterns for Titanic Challange. Below are some links that I used:

"Titanic - Advanced Feature Engineering Tutorial"

The attributes SibSp and Parch represent the count of siblings, spouses, parents, and children that were accompanying the passengers during the voyage. We can merge these attributes to create a new feature.

```
1 data['Age'].fillna(data['Age'].mean(), inplace=True)
2 data['Embarked'].fillna(data['Embarked'].mode()[0], inplace=True)
3 test['Age'].fillna(test['Age'].mean(), inplace=True)
4 test['Embarked'].fillna(test['Embarked'].mode()[0], inplace=True)
```

Here is Kaggle.com scoreboard progress after this implementation.



3 Different Machine Learning Models

3.1 GradientBoostingClassifier

Added this code section then calculated the score on Kaggle.com

```
1 from sklearn.ensemble import GradientBoostingClassifier
2
3 model = GradientBoostingClassifier(n_estimators=100, max_depth=5, random_state=1)
```

Score: 0.76315

[&]quot;Kaggle-Competition-Titanic Survival"

3.2 Support Vector Machine

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.pipeline import make_pipeline
3 from sklearn.svm import SVC
4
5 model = make_pipeline(StandardScaler(), SVC(kernel='rbf', C=1, gamma=0.1, random_state=1)
```

Score: 0.78468

3.3 Naive Bayes

```
1 from sklearn.naive_bayes import GaussianNB
2
3 model = GaussianNB()
```

Score: 0.76076

3.4 Logistic Regression

```
1 from sklearn.linear_model import LogisticRegression
2
3 model = LogisticRegression()
```

Score: 0.75358

4 Hyperparameter Tuning for the Best Scored ML Method

Added n_jobs = -1 parameter to RandomForest model. It's mean "using all processors". I did not found more useful parameters for RandomForest. Then, I Wrote a code for best parameter selection.

```
Parameter Selector
  param_grid = {
       'n_estimators': [100, 300, 500],
       'max_depth': [5, 10, 15],
       'min_samples_split': [2, 5, 10],
       'min_samples_leaf': [1, 2, 4],
       'n_{jobs}': [-1],
       'random_state':[100, 500, 5000]
10
  }
11
  model = RandomForestClassifier(random_state=42)
13
  grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5)
15
  grid_search.fit(X, y)
17
  print("Best parameters: ", grid_search.best_params_)
  print("Best score: ", grid_search.best_score_))
```

After, these parameter tunings we reach the **0.78947** score.

5 Final Python Code

```
#!/usr/bin/env python3
  # @Author: Huseyin Tutan 110190021
  # IMPORT THE LIBRARIES
7 import pandas as pd
8 import numpy as np
 import matplotlib.pyplot as plt
10 import colorama
11
12 from sklearn.metrics import accuracy_score
13 from sklearn.model_selection import cross_val_score
14 from colorama import Fore
15 from sklearn.ensemble import RandomForestClassifier
16 from yellowbrick.model_selection import CVScores
17 from sklearn.model_selection import StratifiedKFold
18 from sklearn.ensemble import GradientBoostingClassifier
19 from sklearn.preprocessing import StandardScaler
20 from sklearn.pipeline import make_pipeline
21 from sklearn.svm import SVC
22 from sklearn.naive_bayes import GaussianNB
23 from sklearn.linear_model import LogisticRegression
24 from sklearn.model_selection import GridSearchCV
26
  # IMPORT THE DATA
  data = pd.read_csv("train.csv")
  test = pd.read_csv("test.csv")
  # MISS VALUES
31
  data['Age'].fillna(data['Age'].mean(), inplace=True)
  test['Age'].fillna(test['Age'].mean(), inplace=True)
  data['Embarked'].fillna(data['Embarked'].mode()[0], inplace=True)
  test['Embarked'].fillna(test['Embarked'].mode()[0], inplace=True)
38
39
```

```
40 # CODE SECTION
 data2 = [data, test]
  for dataset in data2:
       dataset['relatives'] = dataset['SibSp'] + dataset['Parch']
       dataset.loc[dataset['relatives'] > 0, 'travelled_alone'] = 'No'
       dataset.loc[dataset['relatives'] == 0, 'travelled_alone'] = 'Yes'
46
47
48 y = data["Survived"]
49 features = ["Pclass", "Sex", "Age", "Embarked", "travelled_alone"]
  X = pd.get_dummies(data[features])
  X_test = pd.get_dummies(test[features])
  # PARAMETER SELECTOR
53
54
  # param_grid = {
55
         'n_estimators': [100, 300, 500],
56
         'max_depth': [3, 4, 5],
57
         'min_samples_split': [2, 5, 10],
58
         'min_samples_leaf': [1, 2, 4],
         'n_jobs': [-1],
         'random_state':[0, 100, 500, 5000]
62
  # }
63
  # model = RandomForestClassifier(random_state=42)
64
65
  # grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5)
66
67
   # grid_search.fit(X, y)
69
   # print("Best parameters: ", grid_search.best_params_)
   # print("Best score: ", grid_search.best_score_)
72
  model = RandomForestClassifier(n_estimators=100,
                                   max_depth=4,
74
                                   min_samples_split=10,
75
                                   min_samples_leaf=2,
76
                                   random_state=5000,
77
                                   n_{\text{jobs}}=-1)
  # model = GradientBoostingClassifier(n_estimators=100, max_depth=5, random_state=1)
  # model = make_pipeline(StandardScaler(), SVC(kernel='rbf', C=1, gamma=0.1, random_state
  # model = GaussianNB()
82 # model = LogisticRegression()
```

```
83
84
85 model.fit(X, y)
  predictions = model.predict(X_test)
   output = pd.DataFrame(
       {'PassengerId': test.PassengerId, 'Survived': predictions})
   output.to_csv('submission1.csv', index=False)
90
   # COMPUTE CROSS-VALIDATION SCORE
91
92
  scores = cross_val_score(model, X, y, cv=5)
94 print("Cross-validation scores:", scores)
   print("Mean score:", scores.mean())
   # COMPUTE ACCURACY SCORE
98
99 train_predictions = model.predict(X)
100 accuracy = accuracy_score(y, train_predictions)
101 print("Train accuracy:", accuracy)
```