

I decided to develop a classification model for 'cses4_cut.csv' data. First, I constructed features and target matrices. In features matrix, all columns except 'age' have categorical information. Therefore, I used one-hot encoding on features matrix.

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.metrics import accuracy_score
5 import seaborn as sns
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.preprocessing import OneHotEncoder
8
9 #read from csv file
10 df = pd.read_csv('cses4_cut.csv')
11 #construct features matrix, my features are D2003: education and age
12 X = df.iloc[:,1:32]
13 #target matrix
14 Y = df['voted']
15
16 # education is categorical, therefore I will apply one-hot encoding
17 # creating instance of one-hot-encoder
18 X_encoded = OneHotEncoder(sparse=True, handle_unknown='ignore')
19 # passing bridge-types-cat column
20 X_encoded = pd.DataFrame(X_encoded.fit_transform(X.iloc[:,1:31]).toarray())
21 # merge with main df bridge_df on key values
22 X_encoded = X_encoded.join(X['age'])
```

Then, I decided to use Gaussian naive Bayes model as classifier. I split existing data as train (default %75) data and test (default %25) data. I fitted the model, predicted values for Xtest, calculated accuracy of model. Accuracy score was low (appr. 0.32).

```
1 #First, I will use Gaussian naive Bayes model for classification
2
3 Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_encoded, Y,
4 random_state=1) #split the data as train and test
5
6 model = GaussianNB() #used Gaussian naive Bayes model for classification
7 model.fit(Xtrain, Ytrain) #fit the model
8 Y_model = model.predict(Xtest) #predict y values for x test
9
10 accuracy_score(Ytest, Y_model) #calculates accuracy score of mode
```

0.31769996787664634

Then, I decided to use KNeighborsClassifier as classifier. I split existing data as train (default %75) data and test (default %25) data. I fitted the model, predicted values for Xtest, calculated accuracy of model. Accuracy score was high (appr. 0.8).

```
1 #Second, I will use KNeighborsClassifier model for classification
2
3 X_KNtrain, X_KNtest, Y_KNtrain, Y_KNtest = train_test_split(X_encoded, Y,
4 random_state=1)
5 model2 = KNeighborsClassifier(n_neighbors=1)
6 model2.fit(X_KNtrain, Y_KNtrain)
7
8 Y_KN_model = model2.predict(X_KNtest)
9 accuracy_score(Y_KNtest, Y_KN_model)
```

0.799229039511725

I decided to optimize KNeighborsClassifier model because it has high accuracy score. For optimization, I used GridSearchCV module. By help of this module, I tried following hyperparameters:

- **'n_neighbors'**: np.arange(20),
- **'weights'**: ['uniform', 'distance'],
- **'algorithm'** : ['auto', 'ball_tree', 'kd_tree', 'brute']

GridSearchCV module returned best parameters as following:

- **'algorithm'**: 'auto',
- **'n_neighbors'**: 18,
- **'weights'**: 'distance'

```

1  #For model optimization, I will use GridSearchCV
2
3  from sklearn.model_selection import GridSearchCV
4  import numpy as np
5
6  grid_params = {'n_neighbors': np.arange(20),
7                 'weights': ['uniform', 'distance'],
8                 'algorithm': ['auto', 'ball_tree', 'kd_tree', 'brute']}
9
10 grid = GridSearchCV(KNeighborsClassifier(), grid_params, verbose = 1, cv = 7, n_jobs = -1)
11
12 grid.fit(X_KNtrain, Y_KNtrain)
13 grid.best_params_

```

Fitting 7 folds for each of 160 candidates, totalling 1120 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 34 tasks      | elapsed: 11.3s
[Parallel(n_jobs=-1)]: Done 184 tasks    | elapsed: 1.5min
[Parallel(n_jobs=-1)]: Done 434 tasks    | elapsed: 3.8min
[Parallel(n_jobs=-1)]: Done 784 tasks    | elapsed: 7.1min
[Parallel(n_jobs=-1)]: Done 1120 out of 1120 | elapsed: 9.2min finished

```

```
{'algorithm': 'auto', 'n_neighbors': 18, 'weights': 'distance'}
```

Finally, I run KNeighborsClassifier with best hyperparameter values. I got better accuracy score (appr. 0.86).

```

1  #Third, I will use KNeighborsClassifier model with best parameters
2
3  X_KNtrain2, X_KNtest2, Y_KNtrain2, Y_KNtest2 = train_test_split(X_encoded, Y,
4                          random_state=1)
5  model = KNeighborsClassifier(n_neighbors=18, weights='distance', algorithm='auto')
6  model.fit(X_KNtrain2, Y_KNtrain2)
7
8  Y_KN_model2 = model.predict(X_KNtest2)
9  accuracy_score(Y_KNtest2, Y_KN_model2)

```

```
0.8586572438162544
```