

I decided to develop a classification model for 'cses4\_cut.csv' data. My input variables are 'D2003' (EDUCATION), 'D2010' (CURRENT EMPLOYMENT STATUS), 'age'. My target variable is 'voted'.

First, I replaced missing values in 'D2003', 'D2010' columns with most frequent items in each column. Then I constructed features and target matrices. In features matrix, 'D2003', 'D2010' columns 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 from sklearn.impute import SimpleImputer
9
10 #read from csv file
11 df = pd.read_csv('cses4_cut.csv')
12 #construct features matrix, my features are D2003: education and age
13 X = df[['D2003', 'D2010', 'age']]
14 #target matrix
15 Y = df['voted']
16
17 #for D2003 and D2010 columns replaced missing values with most_frequent items
18 X_processed = SimpleImputer(strategy='most_frequent', missing_values=99)
19 X_processed = X_processed.fit(X[['D2003']])
20 X[['D2003']] = X_processed.transform(X[['D2003']])
21 X_processed = X_processed.fit(X[['D2010']])
22 X[['D2010']] = X_processed.transform(X[['D2010']])
23
24 # education is categorical, therefore I will apply one-hot encoding
25 # creating instance of one-hot-encoder
26 X_encoded = OneHotEncoder(sparse=True, handle_unknown='ignore')
27 # passing bridge-types-cat column
28 X_encoded = pd.DataFrame(X_encoded.fit_transform(X[['D2003', 'D2010']]).toarray())
29 # merge encoded columns with age column
30 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 appr. 0.79.

```
1 #split the data as train and test
2 Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_encoded, Y,
3 random_state=1)
```

```
1 #First, I will use Gaussian naive Bayes model for classification
2 model = GaussianNB() #used Gaussian naive Bayes model for classification
3 model.fit(Xtrain, Ytrain) #fit the model
4 Y_model = model.predict(Xtest) #predict y values for x test
5
6 accuracy_score(Ytest, Y_model) #calculates accuracy score of mode
```

0.7895920334082879

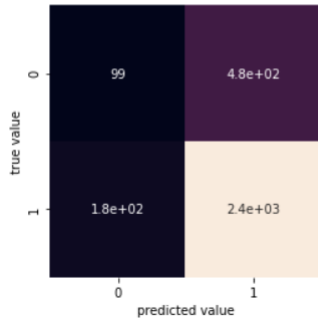
I plotted confusion matrix to see model's output.

```

1 #where did we go wrong? Confusion matrix shows frequency of misclassification
2 import seaborn as sns
3 from sklearn.metrics import confusion_matrix
4 mat = confusion_matrix(Ytest, Y_model)
5 sns.heatmap(mat, square=True, annot=True, cbar=False)
6 plt.xlabel('predicted value')
7 plt.ylabel('true value')

```

Text(91.68, 0.5, 'true value')



Then, I decided to use KNeighborsClassifier as classifier. I fitted the model, predicted values for Xtest, calculated accuracy of model. Accuracy score was appr. 0.77.

```

1 #Second, I will use KNeighborsClassifier model for classification
2
3 model2 = KNeighborsClassifier(n_neighbors=1)
4 model2.fit(Xtrain, Ytrain)
5
6 Y_KN_model = model2.predict(Xtest)
7 accuracy_score(Ytest, Y_KN_model)

```

0.7654995181496949

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': 'brute',
- 'n\_neighbors': 17,
- 'weights': 'uniform'

```

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(Xtrain, Ytrain)
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: 2.1s
[Parallel(n_jobs=-1)]: Done 560 tasks | elapsed: 6.3s
[Parallel(n_jobs=-1)]: Done 1120 out of 1120 | elapsed: 20.2s finished

```

```
{'algorithm': 'brute', 'n_neighbors': 17, 'weights': 'uniform'}
```

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

```

1 #Finally, I will use KNeighborsClassifier model with best parameters
2
3 model = KNeighborsClassifier(n_neighbors=17, weights='uniform', algorithm='brute')
4 model.fit(Xtrain, Ytrain)
5
6 Y_KN_model2 = model.predict(Xtest)
7 accuracy_score(Ytest, Y_KN_model2)

```

```
0.8445229681978799
```

Then I plotted confusion matrix to see model's output.

```

1 #where did we go wrong? Confusion matrix shows frequency of misclassification
2 mat = confusion_matrix(Ytest, Y_KN_model2)
3 sns.heatmap(mat, square=True, annot=True, cbar=False)
4 plt.xlabel('predicted value')
5 plt.ylabel('true value')

```

```
Text(91.68, 0.5, 'true value')
```

