```
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5
6
7
8   from sklearn.model_selection import KFold,cross_val_score, RepeatedStratifiedKFold,StratifiedK
9   from sklearn.impute import SimpleImputer
10  from sklearn.pipeline import Pipeline
11  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
12  from sklearn.preprocessing import OneHotEncoder,StandardScaler,PowerTransformer
13  from sklearn.compose import ColumnTransformer
14  from sklearn.pipeline import Pipeline
15  from sklearn.linear_model import LogisticRegression
16  from sklearn.svm import SVC
17  from sklearn.impute import SimpleImputer
18  from sklearn.dummy import DummyClassifier
19  from imblearn.over_sampling import SMOTE
20
21  from sklearn.ensemble import AdaBoostClassifier
22  from sklearn.ensemble import GradientBoostingClassifier
23  from sklearn.ensemble import RandomForestClassifier
24  from sklearn.ensemble import ExtraTreesClassifier
25  from sklearn.neighbors import KNeighborsClassifier
26
27  import optuna
28  from xgboost import XGBClassifier
29  from lightgbm import LGBMClassifier
30  from catboost import CatBoostClassifier
31
32  from sklearn.pipeline import make_pipeline
33  from sklearn.pipeline import Pipeline
34  from sklearn.compose import make_column_transformer
35
36  from sklearn.model_selection import KFold, cross_val_predict, train_test_split,GridSearchCV,cr
37  from sklearn.metrics import accuracy_score,classification_report
38
39  #importing plotly and cufflinks in offline mode
40  import cufflinks as cf
41  import plotly.offline
42  cf.go_offline()
43  cf.set_config_file(offline=False, world_readable=True)
44
45
46  import plotly
47  import plotly.express as px
48  import plotly.graph_objs as go
49  import plotly.offline as py
50  from plotly.offline import iplot
51  from plotly.subplots import make_subplots
52  import plotly.figure_factory as ff
53
54  import missingno as msno
```

```
54    import missingno as msno
55
56    import warnings
57    warnings.filterwarnings("ignore")
```

⊏→

```
1    pd.set_option('max_columns',100)
2    pd.set_option('max_rows',900)
3
4    pd.set_option('max_colwidth',200)
5
6    df = pd.read_csv('heart.csv')
7    df.head()
```

```
1    df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 918 entries, 0 to 917
Data columns (total 12 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   Age            918 non-null    int64
 1   Sex            918 non-null    object
 2   ChestPainType  918 non-null    object
 3   RestingBP      918 non-null    int64
 4   Cholesterol    918 non-null    int64
 5   FastingBS      918 non-null    int64
 6   RestingECG     918 non-null    object
 7   MaxHR          918 non-null    int64
 8   ExerciseAngina 918 non-null    object
 9   Oldpeak        918 non-null    float64
 10  ST_Slope       918 non-null    object
 11  HeartDisease   918 non-null    int64
dtypes: float64(1), int64(6), object(5)
memory usage: 86.2+ KB
```

```
1 df.duplicated().sum()
```

```
0
```

```
1 def missing (df):
2     missing_number = df.isnull().sum().sort_values(ascending=False)
3     missing_percent = (df.isnull().sum()/df.isnull().count()).sort_values(ascending=False)
```

```
4      missing_values = pd.concat([missing_number, missing_percent], axis=1, keys=['Missing_Number
5      return missing_values
6
7 missing(df)
```

```
1 numerical= df.drop(['HeartDisease'], axis=1).select_dtypes('number').columns
2
3 categorical = df.select_dtypes('object').columns
4
5 print(f'Numerical Columns:  {df[numerical].columns}')
6 print('\n')
7 print(f'Categorical Columns: {df[categorical].columns}')
```

```
Numerical Columns:  Index(['Age', 'RestingBP', 'Cholesterol', 'FastingBS', 'MaxHR', 'Oldpeak'

Categorical Columns: Index(['Sex', 'ChestPainType', 'RestingECG', 'ExerciseAngina', 'ST_Slope
```

```
1 df[categorical].nunique()
```

```
Sex               2
ChestPainType     4
RestingECG        3
ExerciseAngina    2
ST_Slope          3
dtype: int64
```

```
1 y = df['HeartDisease']
2 print(f'Percentage of patient had a HeartDisease:  {round(y.value_counts(normalize=True)[1]*100
```

```
Percentage of patient had a HeartDisease:  55.34 %  --> (508 patient)
Percentage of patient did not have a HeartDisease: 44.66  %  --> (410 patient)
```

```
1 df['HeartDisease'].iplot(kind='hist')
```

```
1  df[numerical].describe()
```

```
1 df['HeartDisease'].iplot(kind='hist')
```

```
1  df[numerical].iplot(kind='hist');
```

```
1  df[numerical].iplot(kind='histogram',subplots=True,bins=50)
```

```
1 skew_limit = 0.75 # This is our threshold-limit to evaluate skewness. Overall below abs(1) seer
2 skew_vals = df[numerical].drop('FastingBS', axis=1).skew()
3 skew_cols= skew_vals[abs(skew_vals)> skew_limit].sort_values(ascending=False)
4 skew_cols
```

```
Oldpeak     1.022872
dtype: float64
```

```
1   numerical1= df.select_dtypes('number').columns
2
3
4   matrix = np.triu(df[numerical1].corr())
5   fig, ax = plt.subplots(figsize=(14,10))
6   sns.heatmap (df[numerical1].corr(), annot=True, fmt= '.2f', vmin=-1, vmax=1, center=0, cmap='c
```

```
1 df[categorical].head()
```

## Gender and Heart Disease

```
1  print (f'A female person has a probability of {round(df[df["Sex"]=="F"]["HeartDisease"].mean()
2
3  print()
4
5  print (f'A male person has a probability of {round(df[df["Sex"]=="M"]["HeartDisease"].mean()*1
6
7  print()
8
```

```
A female person has a probability of 25.91 % have a HeartDisease

A male person has a probability of 63.17 % have a HeartDisease
```

```
1  fig = px.histogram(df, x="Sex", color="HeartDisease",width=400, height=400)
2  fig.show()
```

## ▾ Chest Pain Type and Heart Disease

```
1   df.groupby('ChestPainType')['HeartDisease'].mean().sort_values(ascending=False)
```

```
ChestPainType
ASY    0.790323
TA     0.434783
NAP    0.354680
ATA    0.138728
Name: HeartDisease, dtype: float64
```

```
1   fig = px.histogram(df, x="ChestPainType", color="HeartDisease",width=400, height=400)
2   fig.show()
```

## ▾ RestingECG and Heart Disease

```
1   df.groupby('RestingECG')['HeartDisease'].mean().sort_values(ascending=False)
```

```
RestingECG
ST       0.657303
LVH      0.563830
Normal   0.516304
Name: HeartDisease, dtype: float64
```

```
1   fig = px.histogram(df, x="RestingECG", color="HeartDisease",width=400, height=400)
2   fig.show()
```

## ▾ ExerciseAngina and Heart Disease

```
1 df.groupby('ExerciseAngina')['HeartDisease'].mean().sort_values(ascending=False)
```

```
ExerciseAngina
Y    0.851752
N    0.351005
Name: HeartDisease, dtype: float64
```

```
1 fig = px.histogram(df, x="ExerciseAngina", color="HeartDisease",width=400, height=400)
2 fig.show()
```

## ▾ ST_Slope and Heart Disease

```
1  df.groupby('ST_Slope')['HeartDisease'].mean().sort_values(ascending=False)
```

```
ST_Slope
Flat    0.828261
Down    0.777778
Up      0.197468
Name: HeartDisease, dtype: float64
```

```
1  fig = px.histogram(df, x="ST_Slope", color="HeartDisease",width=400, height=400)
2  fig.show()
```

## ▾ Overall Insights from the Exploratory Data Analysis

```
1   accuracy =[]
2   model_names =[]
3
4
5   X= df.drop('HeartDisease', axis=1)
6   y= df['HeartDisease']
7   X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
8
9   ohe= OneHotEncoder()
10  ct= make_column_transformer((ohe,categorical),remainder='passthrough')
11
```

```
12
13    model = DummyClassifier(strategy='constant', constant=1)
14    pipe = make_pipeline(ct, model)
15    pipe.fit(X_train, y_train)
16    y_pred = pipe.predict(X_test)
17    accuracy.append(round(accuracy_score(y_test, y_pred),4))
18    print (f'model : {model} and  accuracy score is : {round(accuracy_score(y_test, y_pred),4)}')
19
20    model_names = ['DummyClassifier']
21    dummy_result_df = pd.DataFrame({'Accuracy':accuracy}, index=model_names)
22    dummy_result_df
```

```
1     accuracy =[]
2     model_names =[]
3
4
5     X= df.drop('HeartDisease', axis=1)
6     y= df['HeartDisease']
7     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
8
9     ohe= OneHotEncoder()
10    ct= make_column_transformer((ohe,categorical),remainder='passthrough')
11
12
13    lr = LogisticRegression(solver='liblinear')
14    lda= LinearDiscriminantAnalysis()
15    svm = SVC(gamma='scale')
16    knn = KNeighborsClassifier()
17
18    models = [lr,lda,svm,knn]
19
20    for model in models:
21        pipe = make_pipeline(ct, model)
22        pipe.fit(X_train, y_train)
23        y_pred = pipe.predict(X_test)
24        accuracy.append(round(accuracy_score(y_test, y_pred),4))
25        print (f'model : {model} and  accuracy score is : {round(accuracy_score(y_test, y_pred),4)
26
27    model_names = ['Logistic','LinearDiscriminant','SVM','KNeighbors']
28    result_df1 = pd.DataFrame({'Accuracy':accuracy}, index=model_names)
29    result_df1
```

```python
1   accuracy =[]
2   model_names =[]
3
4
5   X= df.drop('HeartDisease', axis=1)
6   y= df['HeartDisease']
7   X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
8
9   ohe= OneHotEncoder()
10  s= StandardScaler()
11  ct1= make_column_transformer((ohe,categorical),(s,numerical))
12
13
14  lr = LogisticRegression(solver='liblinear')
15  lda= LinearDiscriminantAnalysis()
16  svm = SVC(gamma='scale')
17  knn = KNeighborsClassifier()
18
19  models = [lr,lda,svm,knn]
20
21  for model in models:
22      pipe = make_pipeline(ct1, model)
23      pipe.fit(X_train, y_train)
24      y_pred = pipe.predict(X_test)
25      accuracy.append(round(accuracy_score(y_test, y_pred),4))
26      print (f'model : {model} and  accuracy score is : {round(accuracy_score(y_test, y_pred),4)
27
28  model_names = ['Logistic_scl','LinearDiscriminant_scl','SVM_scl','KNeighbors_scl']
29  result_df2 = pd.DataFrame({'Accuracy':accuracy}, index=model_names)
30  result_df2
```

```python
1   accuracy =[]
2   model_names =[]
3
4
```

```python
 5    X= df.drop('HeartDisease', axis=1)
 6    y= df['HeartDisease']
 7    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
 8
 9    ohe= OneHotEncoder()
10    ct= make_column_transformer((ohe,categorical),remainder='passthrough')
11
12    ada = AdaBoostClassifier(random_state=0)
13    gb = GradientBoostingClassifier(random_state=0)
14    rf = RandomForestClassifier(random_state=0)
15    et=  ExtraTreesClassifier(random_state=0)
16
17
18
19    models = [ada,gb,rf,et]
20
21    for model in models:
22        pipe = make_pipeline(ct, model)
23        pipe.fit(X_train, y_train)
24        y_pred = pipe.predict(X_test)
25        accuracy.append(round(accuracy_score(y_test, y_pred),4))
26        print (f'model : {model} and  accuracy score is : {round(accuracy_score(y_test, y_pred),4)
27
28    model_names = ['Ada','Gradient','Random','ExtraTree']
29    result_df3 = pd.DataFrame({'Accuracy':accuracy}, index=model_names)
30    result_df3
```

```python
 1    accuracy =[]
 2    model_names =[]
 3
 4
 5    X= df.drop('HeartDisease', axis=1)
 6    y= df['HeartDisease']
 7    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
 8
 9    ohe= OneHotEncoder()
10    ct= make_column_transformer((ohe,categorical),remainder='passthrough')
11
12    xgbc = XGBClassifier(random_state=0)
13    lgbmc=LGBMClassifier(random_state=0)
14
```

```
15
16   models = [xgbc,lgbmc]
17
18   for model in models:
19       pipe = make_pipeline(ct, model)
20       pipe.fit(X_train, y_train)
21       y_pred = pipe.predict(X_test)
22       accuracy.append(round(accuracy_score(y_test, y_pred),4))
23
24   model_names = ['XGBoost','LightGBM']
25   result_df4 = pd.DataFrame({'Accuracy':accuracy}, index=model_names)
26   result_df4
```

```
1    accuracy =[]
2    model_names =[]
3
4
5    X= df.drop('HeartDisease', axis=1)
6    y= df['HeartDisease']
7    categorical_features_indices = np.where(X.dtypes != np.float)[0]
8
9    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
10
11   model = CatBoostClassifier(verbose=False,random_state=0)
12
13   model.fit(X_train, y_train,cat_features=categorical_features_indices,eval_set=(X_test, y_test)
14   y_pred = model.predict(X_test)
15   accuracy.append(round(accuracy_score(y_test, y_pred),4))
16
17   model_names = ['Catboost_default']
18   result_df5 = pd.DataFrame({'Accuracy':accuracy}, index=model_names)
19   result_df5
20
21
```

```
1    def objective(trial):
2        X= df.drop('HeartDisease', axis=1)
3        y= df['HeartDisease']
4        categorical_features_indices = np.where(X.dtypes != np.float)[0]
5
6        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
7
8        param = {
```

```python
 9            "objective": trial.suggest_categorical("objective", ["Logloss", "CrossEntropy"]),
10            "colsample_bylevel": trial.suggest_float("colsample_bylevel", 0.01, 0.1),
11            "depth": trial.suggest_int("depth", 1, 12),
12            "boosting_type": trial.suggest_categorical("boosting_type", ["Ordered", "Plain"]),
13            "bootstrap_type": trial.suggest_categorical(
14                "bootstrap_type", ["Bayesian", "Bernoulli", "MVS"]
15            ),
16            "used_ram_limit": "3gb",
17        }
18
19        if param["bootstrap_type"] == "Bayesian":
20            param["bagging_temperature"] = trial.suggest_float("bagging_temperature", 0, 10)
21        elif param["bootstrap_type"] == "Bernoulli":
22            param["subsample"] = trial.suggest_float("subsample", 0.1, 1)
23
24        cat_cls = CatBoostClassifier(**param)
25
26        cat_cls.fit(X_train, y_train, eval_set=[(X_test, y_test)], cat_features=categorical_featur
27
28        preds = cat_cls.predict(X_test)
29        pred_labels = np.rint(preds)
30        accuracy = accuracy_score(y_test, pred_labels)
31        return accuracy
32
33
34    if __name__ == "__main__":
35        study = optuna.create_study(direction="maximize")
36        study.optimize(objective, n_trials=50, timeout=600)
37
38        print("Number of finished trials: {}".format(len(study.trials)))
39
40        print("Best trial:")
41        trial = study.best_trial
42
43        print("  Value: {}".format(trial.value))
44
45        print("  Params: ")
46        for key, value in trial.params.items():
47            print("    {}: {}".format(key, value))
```

```
[I 2022-04-04 11:09:46,497] Trial 3 finished with value: 0.8659420289855072 and parameter:
[I 2022-04-04 11:09:49,333] Trial 4 finished with value: 0.8804347826086957 and parameter:
[I 2022-04-04 11:09:55,124] Trial 5 finished with value: 0.8913043478260869 and parameter:
[I 2022-04-04 11:09:56,384] Trial 6 finished with value: 0.8876811594202898 and parameter:
[I 2022-04-04 11:09:59,400] Trial 7 finished with value: 0.8876811594202898 and parameter:
[I 2022-04-04 11:10:01,216] Trial 8 finished with value: 0.8876811594202898 and parameter:
[I 2022-04-04 11:10:02,627] Trial 9 finished with value: 0.8913043478260869 and parameter:
[I 2022-04-04 11:10:06,275] Trial 10 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:07,608] Trial 11 finished with value: 0.894927536231884 and parameter:
[I 2022-04-04 11:10:08,988] Trial 12 finished with value: 0.8840579710144928 and paramete
[I 2022-04-04 11:10:10,432] Trial 13 finished with value: 0.8876811594202898 and paramete
[I 2022-04-04 11:10:11,768] Trial 14 finished with value: 0.8876811594202898 and paramete

[I 2022-04-04 11:10:13,133] Trial 15 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:14,787] Trial 16 finished with value: 0.8876811594202898 and paramete
[I 2022-04-04 11:10:16,118] Trial 17 finished with value: 0.8840579710144928 and paramete
[I 2022-04-04 11:10:17,620] Trial 18 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:20,256] Trial 19 finished with value: 0.8768115942028986 and paramete
```

```
[I 2022-04-04 11:10:21,998] Trial 20 finished with value: 0.8804347826086957 and paramete
[I 2022-04-04 11:10:23,460] Trial 21 finished with value: 0.8840579710144928 and paramete
[I 2022-04-04 11:10:24,807] Trial 22 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:26,146] Trial 23 finished with value: 0.8876811594202898 and paramete
[I 2022-04-04 11:10:27,444] Trial 24 finished with value: 0.9021739130434783 and paramete
[I 2022-04-04 11:10:28,674] Trial 25 finished with value: 0.8876811594202898 and paramete
[I 2022-04-04 11:10:29,414] Trial 26 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:34,060] Trial 27 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:35,374] Trial 28 finished with value: 0.8876811594202898 and paramete
[I 2022-04-04 11:10:36,093] Trial 29 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:36,786] Trial 30 finished with value: 0.8840579710144928 and paramete
[I 2022-04-04 11:10:37,530] Trial 31 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:38,303] Trial 32 finished with value: 0.8985507246376812 and paramete
[I 2022-04-04 11:10:39,083] Trial 33 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:39,760] Trial 34 finished with value: 0.8768115942028986 and paramete
[I 2022-04-04 11:10:41,120] Trial 35 finished with value: 0.8804347826086957 and paramete
[I 2022-04-04 11:10:41,925] Trial 36 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:43,250] Trial 37 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:44,048] Trial 38 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:44,767] Trial 39 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:45,545] Trial 40 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:46,383] Trial 41 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:47,223] Trial 42 finished with value: 0.8985507246376812 and paramete
[I 2022-04-04 11:10:48,170] Trial 43 finished with value: 0.8876811594202898 and paramete
[I 2022-04-04 11:10:48,920] Trial 44 finished with value: 0.8804347826086957 and paramete
[I 2022-04-04 11:10:49,806] Trial 45 finished with value: 0.8913043478260869 and paramete
[I 2022-04-04 11:10:50,563] Trial 46 finished with value: 0.8985507246376812 and paramete
[I 2022-04-04 11:10:51,341] Trial 47 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:52,220] Trial 48 finished with value: 0.894927536231884 and parameters
[I 2022-04-04 11:10:53,107] Trial 49 finished with value: 0.8876811594202898 and paramete
Number of finished trials: 50
Best trial:
  Value: 0.9021739130434783
  Params:
    objective: Logloss
    colsample_bylevel: 0.01134943888867012
    depth: 9
    boosting_type: Ordered
    bootstrap_type: MVS
```

```
1 accuracy =[]
2 model_names =[]
3
4
5 X= df.drop('HeartDisease', axis=1)
6 y= df['HeartDisease']
7 categorical_features_indices = np.where(X.dtypes != np.float)[0]
8
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
10
11 model = CatBoostClassifier(verbose=False,random_state=0,
12                     objective= 'CrossEntropy',
13     colsample_bylevel= 0.04292240490294766,
14     depth= 10,
15     boosting_type= 'Plain',
16     bootstrap_type= 'MVS')
```

```
17
18 model.fit(X_train, y_train,cat_features=categorical_features_indices,eval_set=(X_test, y_test)
19 y_pred = model.predict(X_test)
20 accuracy.append(round(accuracy_score(y_test, y_pred),4))
21 print(classification_report(y_test, y_pred))
22
23 model_names = ['Catboost_tuned']
24 result_df6 = pd.DataFrame({'Accuracy':accuracy}, index=model_names)
25 result_df6
26
27
```

```
1 feature_importance = np.array(model.get_feature_importance())
2 features = np.array(X_train.columns)
3 fi={'features':features,'feature_importance':feature_importance}
4 df_fi = pd.DataFrame(fi)
5 df_fi.sort_values(by=['feature_importance'], ascending=True,inplace=True)
6 fig = px.bar(df_fi, x='feature_importance', y='features',title="CatBoost Feature Importance",he
7 fig.show()
```

```
1 result_final = pd.concat([dummy_result_df,result_df1,result_df2,result_df3,result_df4,result_d
```

```
1    result_final.sort_values(by=['Accuracy'], ascending=True,inplace=True)
2    fig = px.bar(result_final, x='Accuracy', y=result_final.index,title='Model Comparison',height=
3    fig.show()
```

✓ 0초  오후 8:10에 완료됨  ● ✕