

Predictive Modeling for Chronic obstructive pulmonary disease Outcomes using Machine Learning

**Problem**

Chronic Obstructive Pulmonary Disease (COPD) is a progressive respiratory condition characterized by persistent airflow limitation, significantly affecting millions worldwide. Despite being preventable and manageable, COPD often remains undiagnosed, especially in its early stages, due to its subtle symptoms and overlapping risk factors such as smoking, environmental pollutants, and genetic factors. Traditional diagnostic approaches, such as spirometry, are often conducted late in the disease progression, limiting opportunities for early intervention. Furthermore, the complex nature of COPD makes its prediction challenging, often leading to delays in appropriate care. There is a critical need for innovative solutions that leverage machine learning to analyze diverse patient data, predict disease onset accurately, and identify individuals at high risk. Developing such models can revolutionize COPD management by facilitating timely interventions, personalized treatment plans, and ultimately improving patient outcomes.

```
!pip install researchpy
```

```
Requirement already satisfied: researchpy in /usr/local/lib/python3.10/dist-packages (0.3.0)
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from researchpy) (1.13.1)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from researchpy) (2.16.4)
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from researchpy) (2.2.3)
Requirement already satisfied: statsmodels in /usr/local/lib/python3.10/dist-packages (from researchpy) (0.14.3)
Requirement already satisfied: pyts in /usr/local/lib/python3.10/dist-packages (from researchpy) (0.15.0)
Requirement already satisfied: pyts-2024.1 in /usr/local/lib/python3.10/dist-packages (from pandas<researchpy> (0.8.2))
Requirement already satisfied: python-dateutil in /usr/local/lib/python3.10/dist-packages (from pandas<researchpy> (0.8.2))
Requirement already satisfied: pytz in /usr/local/lib/python3.10/dist-packages (from pandas<researchpy> (0.8.2))
Requirement already satisfied: tzdata in /usr/local/lib/python3.10/dist-packages (from pandas<researchpy> (0.8.2))
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (from pandas<researchpy> (0.8.2))
Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from statsmodels<researchpy> (0.16.4))
```

```
# Load libraries
import numpy as np
import pandas as pd
import researchpy as rp
import matplotlib.pyplot as plt
import seaborn as sns
```

```
# Load dataset
data = pd.read_csv("content/COPD dataset.csv")
data.head()
```

```
   SL  ID  AGE  PackHistory  COPDSEVERITY  MNT1  MNT2  MNT1Bent  FEV1  FEV1RED  ...  SGQR  AQQuartiles  copd  gender  smoking  Diabetes  muscular  hypertension  AtrialFib  Outcome
0    1  58  77      60.0    SEVERE    120.0  120.0    120.0  1.21  36.0  ...  69.55      4    3    1    2    1    0    0    1    0    1
1    2   79  79      50.0    MODERATE  165.0  176.0    176.0  1.09  56.0  ...  44.24      4    2    0    2    1    0    0    1    1    1
2    3   42  80      11.0    MODERATE  201.0  180.0    201.0  1.32  68.0  ...  44.09      4    2    0    2    1    0    0    1    0    0
3    4  145  96      60.0    VERYSEVERE  210.0  210.0    210.0  0.47  14.0  ...  42.04      1    4    1    2    0    0    0    1    1    0
4    5  136  65      68.0    SEVERE    204.0  210.0    210.0  1.07  42.0  ...  75.56      1    3    1    2    0    1    1    1    0    0
5 rows x 24 columns
```

```
data.drop('SL', axis=1, inplace=True)
```

```
data.drop('ID', axis=1, inplace=True)
```

```
data.head()
```

```
   AGE  PackHistory  COPDSEVERITY  MNT1  MNT2  MNT1Bent  FEV1  FEV1RED  FVC  FVCRED  ...  SGQR  AQQuartiles  copd  gender  smoking  Diabetes  muscular  hypertension  AtrialFib  Outcome
0    77      60.0    SEVERE    120.0  120.0    120.0  1.21  36.0  24.0  98  ...  69.55      4    3    1    2    1    0    0    1    0    1
1    79      50.0    MODERATE  165.0  176.0    176.0  1.09  56.0  1.64  65  ...  44.24      4    2    0    2    1    0    0    1    0    1
2    80      11.0    MODERATE  201.0  180.0    201.0  1.32  68.0  2.30  66  ...  44.09      4    2    0    2    1    0    0    1    0    0
3    16      60.0    VERYSEVERE  210.0  210.0    210.0  0.47  14.0  1.74  27  ...  42.04      1    4    1    2    0    0    0    1    1    0
4    65      68.0    SEVERE    204.0  210.0    210.0  1.07  42.0  2.91  98  ...  75.56      1    3    1    2    0    1    1    1    0    0
5 rows x 22 columns
```

Exploring Data

```
# Check shape of data
data.shape
```

```
(161, 22)
```

```
# dtypes
data.dtypes
```

```
AGE      int64
PackHistory      float64
COPDSEVERITY      object
MNT1      float64
MNT2      float64
MNT1Bent      float64
FEV1      float64
FEV1RED      float64
FVC      float64
FVCRED      int64
CAT      int64
HAD      float64
SGQR      float64
AQQuartiles      int64
copd      int64
gender      int64
smoking      int64
Diabetes      int64
muscular      int64
hypertension      int64
AtrialFib      int64
Outcome      int64
dtype: object
```

```
# info
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 161 entries, 0 to 160
Data columns (total 22 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   AGE         161 non-null    int64
 1   PackHistory 161 non-null    float64
 2   COPDSEVERITY 161 non-null    object
 3   MNT1        161 non-null    float64
 4   MNT2        161 non-null    float64
 5   MNT1Bent    161 non-null    float64
 6   FEV1        161 non-null    float64
 7   FEV1RED     161 non-null    float64
 8   FVC         161 non-null    float64
 9   FVCRED      161 non-null    int64
10   CAT         161 non-null    int64
11   HAD         161 non-null    float64
12   SGQR        161 non-null    float64
13   AQQuartiles 161 non-null    int64
14   copd        161 non-null    int64
15   gender      161 non-null    int64
16   smoking     161 non-null    int64
17   Diabetes    161 non-null    int64
18   muscular    161 non-null    int64
19   hypertension 161 non-null    int64
20   AtrialFib   161 non-null    int64
21   Outcome     161 non-null    int64
dtypes: float64(9), int64(10), object(3)
memory usage: 27.2+ KB
```

```
# Check missing data
data.isnull().sum()
```

```
AGE      0
PackHistory      0
COPDSEVERITY      0
MNT1      0
MNT2      1
MNT1Bent      1
FEV1      0
FEV1RED      0
FVC      0
FVCRED      0
CAT      0
HAD      0
SGQR      0
AQQuartiles      0
copd      0
gender      0
smoking      0
Diabetes      0
muscular      0
hypertension      0
AtrialFib      0
Outcome      0
dtype: int64
```

```
data = data.dropna()
```

```
data[['COPDSEVERITY', 'AQQuartiles', 'copd', 'gender', 'smoking', 'Diabetes', 'muscular', 'hypertension', 'AtrialFib', 'Outcome']] = data[['COPDSEVERITY', 'AQQuartiles', 'copd', 'gender', 'smoking', 'Diabetes', 'muscular', 'hypertension', 'AtrialFib', 'Outcome']]
```

```
<ipython-input-12-478a286dca-1>: Setting up CopyWarning:
  A value is trying to be set on a copy of a slice from a DataFrame.
  Try using .loc[row_index,col_index] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
data[['COPDSEVERITY', 'AQQuartiles', 'copd', 'gender', 'smoking', 'Diabetes', 'muscular', 'hypertension', 'AtrialFib', 'Outcome']] = data[['COPDSEVERITY', 'AQQuartiles', 'copd', 'gender', 'smoking', 'Diabetes', 'muscular', 'hypertension', 'AtrialFib', 'Outcome']].astype(object)
```

```
data.isnull().sum()
```

```
AGE      0
PackHistory      0
COPDSEVERITY      0
MNT1      0
MNT2      0
MNT1Bent      0
FEV1      0
FEV1RED      0
FVC      0
FVCRED      0
CAT      0
HAD      0
SGQR      0
AQQuartiles      0
copd      0
gender      0
smoking      0
Diabetes      0
muscular      0
hypertension      0
AtrialFib      0
Outcome      0
dtype: int64
```

Descriptive Statistic

```
# select numeric data
num_cols = data.select_dtypes(include = "object")
num_cols.head()
```

AGE

PackHistory

MWT1

MWT2

MWTBest

FEV1

FEV1Pred

FVC

FVCpred

CAT

HAD

SGRQ

0	77	60.0	1020.0	1020.0	1.21	36.0	2.40	98	25	8.0	49.50	
1	79	50.0	165.0	176.0	1.09	56.0	1.44	45	12	21.0	44.04	
2	80	11.0	201.0	190.0	201.0	1.52	68.0	2.30	86	22	18.0	44.09
3	56	60.0	210.0	210.0	210.0	0.47	14.0	1.14	27	28	26.0	62.04
4	65	68.0	204.0	210.0	210.0	1.07	42.0	2.81	98	32	18.0	75.56

Next steps: [Generate code with run\\_rtc](#) [View recommended plots](#) [New interactive sheet](#)

run\_cols.columns

Index([AGE, 'PackHistory', 'MWT1', 'MWT2', 'MWTBest', 'FEV1', 'FEV1Pred', 'FVC', 'FVCpred', 'CAT', 'HAD', 'SGRQ'], dtype=object)

# Summary statistics of numerical variables

rp.summary\_statistics(include=['AGE', 'PackHistory', 'MWT1', 'MWT2', 'MWTBest', 'FEV1', 'FEV1Pred', 'FVC', 'FVCpred', 'CAT', 'HAD', 'SGRQ'])

	Name	N	Mean	Median	Variance	SD	SE	95% Conf. Interval	
0	AGE	99	70.1212	71.0	61.8802	7.8721	0.7913		AGE
1	PackHistory	99	39.803	36.0	506.3001	24.6232	2.4747		PackHistory
2	MWT1	99	388.6886	470.0	10071.8471	104.7462	10.2379		MWT1
3	MWT2	99	389.798	399.0	11704.0404	108.1945	10.874		MWT2
4	MWTBest	99	398.7172	420.0	11453.1845	107.0196	10.7359		MWTBest
5	FEV1	99	1.6031	1.6	0.457	0.676	0.0679		FEV1
6	FEV1Pred	99	58.4008	60.0	503.8626	22.4465	2.286		FEV1Pred
7	FVC	99	2.9278	2.77	0.9069	0.970	0.0984		FVC
8	FVCpred	99	86.5434	84.0	481.1053	21.9347	2.2045		FVCpred
9	CAT	99	19.406	18.0	305.5464	18.856	1.8951		CAT
10	HAD	99	11.1191	10.0	73.8956	8.5963	0.864		HAD
11	SGRQ	99	40.1096	38.21	237.075	15.2366	1.5452		SGRQ

# select categorical data

cat\_cols = data.select\_dtypes(include = 'object')

cat\_cols.head()

COPDSEVERITY

AGEquartiles

copd

gender

smoking

diabetes

muscular

hypertension

AtrialFib

Outcome

0	SEVERE	4	3	1	2	1	0	0	1	0
1	MODERATE	4	2	0	2	1	0	0	1	1
2	MODERATE	4	2	0	2	1	0	0	1	0
3	VERY SEVERE	1	4	1	2	0	0	1	1	0
4	SEVERE	1	3	1	2	0	1	1	0	0

Next steps: [Generate code with cat\\_rtc](#) [View recommended plots](#) [New interactive sheet](#)

cat\_cols.columns

Index(['COPDSEVERITY', 'AGEquartiles', 'copd', 'gender', 'smoking', 'diabetes', 'muscular', 'hypertension', 'AtrialFib', 'Outcome'], dtype=object)

# Summary statistics of categorical variables

rp.summary\_statistics(include=['COPDSEVERITY', 'AGEquartiles', 'copd', 'gender', 'smoking', 'diabetes', 'muscular', 'hypertension', 'AtrialFib', 'Outcome'])

	Variable	Outcome	Count	Percent	
0	COPDSEVERITY	MODERATE	41	41.41	
1	SEVERE		27	27.27	
2	MILD		23	23.23	
3	VERY SEVERE		8	8.08	
4	AGEquartiles	3	28	28.28	
5		1	25	25.25	
6		2	24	24.24	
7		4	22	22.22	
8	copd	2	41	41.41	
9		3	27	27.27	
10		1	23	23.23	
11		4	8	8.08	
12	gender	1	63	63.64	
13		0	36	36.36	
14	smoking	2	83	83.84	
15		1	16	16.16	
16	Diabetes	0	78	78.79	
17		1	21	21.21	
18	muscular	0	60	60.61	
19		1	19	19.19	
20	hypertension	0	88	88.89	
21		1	11	11.11	
22	AtrialFib	0	79	79.80	
23		1	20	20.20	
24	Outcome	0	90	90.91	
25		1	9	9.09	

Correlations between Variable

# correlation Pearson's by default

run\_cols.corr(method='pearson')

AGE

PackHistory

MWT1

MWT2

MWTBest

FEV1

FEV1Pred

FVC

FVCpred

CAT

HAD

SGRQ

AGE	1.000000	-0.021614	-0.246847	-0.234862	-0.227397	0.087169	0.080362	-0.129459	0.001331	-0.219833	0.127909
PackHistory	-0.021614	1.000000	-0.251895	-0.280035	-0.249815	-0.132797	-0.123172	-0.091967	0.001033	-0.146967	0.044868
MWT1	-0.246847	-0.251895	1.000000	0.945484	0.902271	0.888609	0.842071	0.498901	-0.110261	-0.262140	0.107047
MWT2	-0.234862	-0.280035	0.945484	1.000000	0.982544	0.471969	0.490252	0.454576	0.291864	-0.151357	-0.293152
MWTBest	-0.227397	-0.249815	0.902271	0.982544	1.000000	0.468608	0.386947	0.444557	0.258128	-0.159393	-0.294838
FEV1	0.087169	-0.132797	0.468605	0.470493	0.468083	1.000000	0.773523	0.0118421	0.514694	-0.063451	0.163453
FEV1Pred	0.080362	-0.123172	0.842057	0.490255	0.389467	0.773523	1.000000	0.519198	0.624727	-0.090657	-0.121885
FVC	-0.129459	-0.091967	0.498909	0.490259	0.444557	0.0118421	0.519199	1.000000	0.624728	-0.140500	0.220218
FVCpred	0.001331	0.001033	0.257159	0.299684	0.259128	0.514694	0.624727	0.623930	1.000000	-0.137885	0.160940
CAT	-0.001331	-0.146967	-0.135281	-0.151937	-0.153993	0.063451	0.090657	-0.188801	-0.137885	1.000000	0.160804
HAD	-0.219833	0.044868	-0.262140	-0.293162	-0.294838	-0.163453	-0.121885	-0.145320	-0.160940	0.160804	1.000000
SGRQ	0.127909	0.042651	0.107541	-0.323944	-0.337424	-0.314719	-0.340204	-0.222182	-0.294018	0.290358	0.395066

Skewness

run\_cols.skew()

AGE

PackHistory

MWT1

MWT2

MWTBest

FEV1

FEV1Pred

FVC

FVCpred

CAT

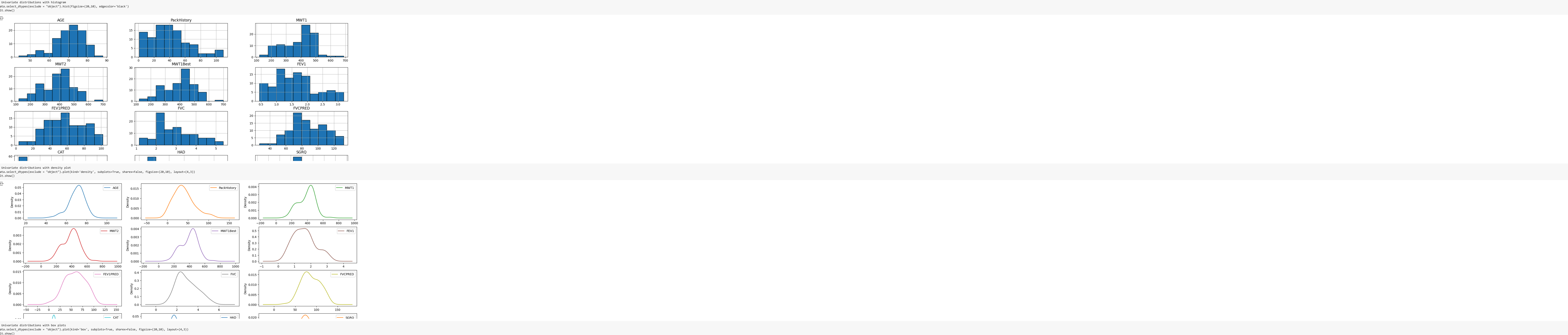
HAD

SGRQ

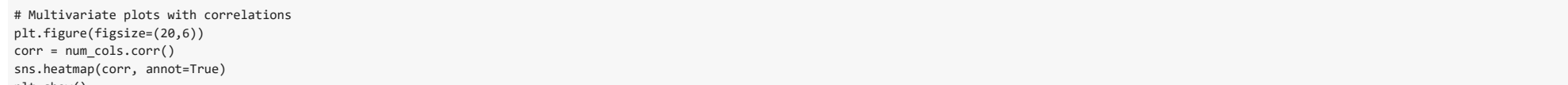
AGE	-0.737041
PackHistory	0.761152
MWT1	-0.311096
MWT2	-0.194043
MWTBest	-0.270568
FEV1	0.463046
FEV1Pred	-0.102857
FVC	0.309557
FVCpred	0.011179
CAT	7.386027
HAD	1.807793
SGRQ	0.191143

dtype: float64

Data Visualizations





[illegible][illegible]

```

$ pip install --upgrade pandas
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (2.2.3)
Requirement already satisfied: numpy>=1.22.4 in /usr/local/lib/python3.10/dist-packages (from pandas) (1.26.4)
Requirement already satisfied: python-dateutil<2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas) (2024.1)
Requirement already satisfied: numpy>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil<2.8.2,>=2.8.2-pandas) (1.16.0)

```

E (compare baseline models)										
Test = compare_model()										
	Model	Accuracy	AUC	Recall	Prec	T1	CPUs	NCE	TT (sec)	
lightgbm	Light Gradient Boosting Machine	0.938	0.933	0.100	0.100	0.100	nan	0.1000	0.0210	
	K-Neighbors Classifier	0.943	0.983	0.000	0.000	0.000	nan	0.0000	0.1500	
dummy	Quadratic Discriminant Analysis	0.943	0.900	0.000	0.000	0.000	nan	0.0000	0.1600	
	Dummy Classifier	0.943	0.900	0.000	0.000	0.000	nan	0.0000	0.2700	
svm	Linear Kernel	0.905	0.567	0.000	0.000	0.000	nan	0.0000	0.1460	
	Ridge Classifier	0.976	0.490	0.000	0.000	0.000	nan	0.0000	0.1500	
ada	Random Forest Classifier	0.976	0.370	0.000	0.000	0.000	nan	0.0000	0.3060	
	Ada Boost Classifier	0.976	0.500	0.100	0.100	0.100	nan	0.1000	0.3020	
sgd	Extra Trees Classifier	0.976	0.424	0.000	0.000	0.000	nan	0.0000	0.2710	
	Exponential Gradient Boosting	0.967	0.400	0.000	0.000	0.000	nan	0.0000	0.1460	
gb	Logistic Regression	0.924	0.933	0.200	0.150	0.167	nan	0.147	0.8130	
	Gradient Boosting Classifier	0.934	0.933	0.000	0.000	0.000	nan	-0.003	0.2960	
lfa	Linear Discriminant Analysis	0.962	0.400	0.100	0.233	0.000	nan	0.000	0.1640	
	Decision Tree Classifier	0.702	0.847	0.000	0.823	0.147	nan	0.040	0.2700	
nb	Naïve Bayes	0.843	0.433	0.400	0.163	0.170	0.058	0.058	0.2400	

Create Model

```
lightgbm = create_model(lightgbm)

Accuracy  AUC Recall  Prec.  F1  Kappa  MCC
Fold
0  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
1  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
2  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
3  0.8571  1.0000  0.0000  0.0000  0.0000  0.0000  0.0000
4  0.8571  0.2523  0.0000  0.0000  0.0000  0.0000  0.0000
5  1.0000  1.0000  1.0000  1.0000  1.0000  1.0000  1.0000
6  0.8571  0.1667  0.0000  0.0000  0.0000  0.0000  0.0000
7  0.8571  0.2523  0.0000  0.0000  0.0000  0.0000  0.0000
8  0.8571  1.0000  0.0000  0.0000  0.0000  0.0000  0.0000
9  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.1000
Mean  0.9286  0.2823  0.1000  0.1000  0.1000  run 0.1000
Std  0.0714  0.4230  0.2000  0.2000  0.3000  run 0.3000

# print model parameters
print(lightgbm)

LGBMClassifier(boosting_type='gbdt', class_weight=None, colsample_bytree=1.0,
               importance_type='split', learning_rate=0.1, max_depth=1,
               min_child_samples=20, min_child_weight=0.001, min_split_gain=0.0,
               n_estimators=100, n_jobs=-1, num_leaves=30, objective=None,
               random_state=123, reg_alpha=0.0, reg_lambda=0.0, subsample=1.0,
               verbose=10, xgb_model=None, xgb_model_type='gbdt')
```

Tune Model

```
tuned_lightgbm = tune_model(lightgbm)

Accuracy  AUC Recall  Prec.  F1  Kappa  MCC
Fold
0  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
1  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
2  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
3  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
4  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
5  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
6  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
7  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
8  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
9  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
Mean  0.9143  0.3000  0.0000  0.0000  0.0000  run 0.0000
Std  0.0702  0.2449  0.0000  0.0000  0.0000  run 0.0000

Fitting 18 folds for each of 18 candidates, totalling 324 fits
Original model was better than the tuned model, hence it will be returned. NOTE: The display metrics are for the tuned model (not the original one).
```

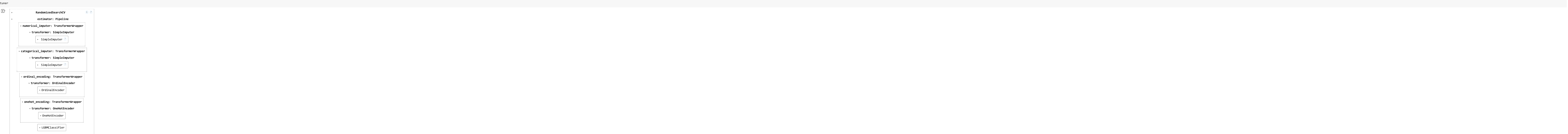
```
tuned_lightgbm_tuner = tune_model(lightgbm, return_tuner=True)

Accuracy  AUC Recall  Prec.  F1  Kappa  MCC
Fold
0  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
1  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
2  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
3  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
4  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
5  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
6  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
7  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
8  0.8571  0.5000  0.0000  0.0000  0.0000  0.0000  0.0000
9  1.0000  0.0000  0.0000  0.0000  0.0000  run 0.0000
Mean  0.9143  0.3000  0.0000  0.0000  0.0000  run 0.0000
Std  0.0702  0.2449  0.0000  0.0000  0.0000  run 0.0000

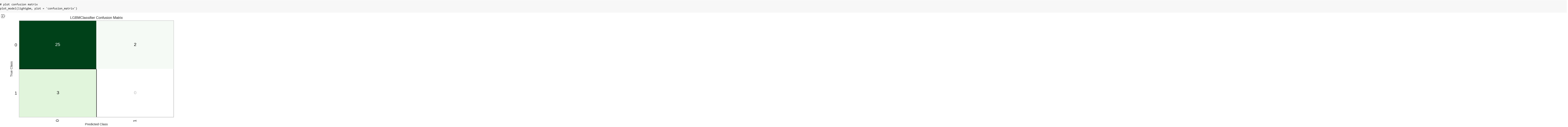
Fitting 18 folds for each of 18 candidates, totalling 324 fits
Original model was better than the tuned model, hence it will be returned. NOTE: The display metrics are for the tuned model (not the original one).
```

```
tuned_lightgbm

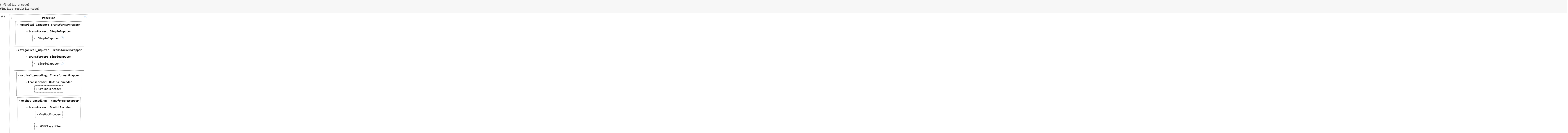
LGBMClassifier
-
LGBMClassifier(boosting_type='gbdt', class_weight=None, colsample_bytree=1.0,
               importance_type='split', learning_rate=0.1, max_depth=1,
               min_child_samples=20, min_child_weight=0.001, min_split_gain=0.0,
               n_estimators=100, n_jobs=-1, num_leaves=30, objective=None,
               random_state=123, reg_alpha=0.0, reg_lambda=0.0, subsample=1.0,
               verbose=10, xgb_model=None, xgb_model_type='gbdt')
```



Analyze Model



Evaluate Model



Prediction

```
# predict on test set
holdout_pred = predict_model(lightgbm)

Model Accuracy  AUC Recall  Prec.  F1  Kappa  MCC
0  Light Gradient Boosting Machine  0.8333  0.4444  0.0000  0.0000  0.0000  0.0070  0.0891

# show predictions
holdout_pred.head()
```



	Ace	Fractional	COPDSEVERITY	MWt1	MWt2	Hemifist	FVEF1	FVEF2	FVC	FVCED	...	copd	gender	smoking	Diabetes	muscular	hypertension	Atrialfibril	Outcome	prediction_Lab1	prediction_score
37	74	640	MODERATE	3840	3720	3840	2700	470	1110	0	0	0	2	1	2	0	0	0	0	0	0.7810
44	73	930	MODERATE	4005	4045	4045	2403	750	1537	125	-	2	1	2	0	0	0	0	0	0	0.7821
76	75	510	MILD	3900	3960	3960	2910	132	1110	0	0	0	1	2	1	0	0	0	0	0	0.6385
92	71	510	MODERATE	4002	4210	4020	2343	790	1433	107	-	2	1	2	0	0	0	0	0	0	0.6904
90	65	340	MODERATE	4200	4230	4200	1445	610	2855	101	-	2	0	1	0	0	0	0	0	0	0.6912

5 rows x 24 columns

```
# copy data and drop Class variable
new_data = data.copy()
new_data.drop('Outcome', axis=1, inplace=True)
new_data.head()
```

AGE	PROCTITIS	CONDISYMPHY	MUT1	MUT2	FEV1	FEV1/FVC	FEV1/FVC2	HA0	SG26	ADG	QUINYSLE	copd	gender	smoking	diabetes	muscular	vascular	atrialfib
77	60.0	SEVERE	120.0	12.0	120.0	1.21	36.0	24.0	NA	80	95.93	4	3	1	2	1	0	1
1	79	50.0	MODERATE	160.0	170.0	170.0	19.0	54.0	116.4	65	-21.0	42.4	4	2	0	0	0	0
2	76	11.0	MODERATE	160.0	190.0	201.0	1.82	60.0	230	96	-18.0	40.0	2	2	2	2	0	0
3	76	50.0	MODERATE	210.0	210.0	210.0	0.87	140.0	114	27	-20.0	42.04	1	2	0	0	0	1
4	65	60.0	SEVERE	250.0	210.0	210.0	0.87	42.0	21.1	98	-18.0	75.56	1	3	1	2	1	1

5 rows + 21 columns

```
# predict model on new_data
predictions = predict_model(lightgbm, data = new_data)
print(predictions.head())
```

AGE	Sex	Acquiescence	COMORBIDITY	MUT1	MUT2	MUT3	FEV1	FVC	FEV1/FVC	FEF25-75	PPFEV1	...	ADQnaires	copd	anger	anxiety	Diabetes	muscular	hypertension	Atrialfibr	stroke	schizophrenia	prediabetic	score
0	77	60.0	MODERATE	120.0	12.0	12.0	12.0	1.21	36.0	24.0	...	...	4	3	2	1	0	0	0	1	0	0	0	0.580
1	79	50.0	MODERATE	145.0	17.0	17.0	17.0	1.60	50.0	36.0	14.4	6.5	...	4	2	0	2	0	0	0	1	0	0	0.9166
2	80	11.0	MODERATE	201.0	19.0	19.0	19.0	1.82	62.0	39.0	22.0	8.6	...	4	2	0	2	0	0	0	1	0	0	0.907
3	84	68.0	SEVERE	215.0	21.0	21.0	21.0	2.03	67	44.0	14.4	11.4	...	2	1	4	1	2	0	0	0	0	0	0.781
4	65	48.0	MODERATE	204.0	21.0	21.0	21.0	2.02	42.0	29.1	96	...	...	1	3	1	2	0	1	1	0	0	0	0.7161

5 rows + 23 columns

## Save Model

```
# save pipeline
save_model(lightgbm, 'E:\COPD_Model')
```

Transformation Pipeline and Model Successfully Saved  
(closing/opening: Memory-Clearing: none)

```
steps=[['numerical_imputer',
        'transformer_encoder', 'overfit_detector']
```

```
include=['AGE', 'PackHistory', 'MWT1',
        'MWT2', 'MWT1Best', 'EPV3']
```

'FEV1PRED', 'FVC', 'FVC/PRED',  
'CAT', 'HMO', 'SZRO']

```
transformer=SimpleInputer(add_indicator=False,  
                           copy=True,
```

```
fill_value=None,
keep_empty_features=False
```

```
missing_values=ran,
strategy='mean'))),
```

```

{catogo...
  UGBClassifier(boosting_type="gbdt", class_weight=None,
    colsample_bytree=0.1, subsample=0.5, min_samples_split=10,

```

```

    (colsample_bytree=1.0, max_depth=5,
      learning_rate=0.1, max_depth=1,
      min_child_weight=10, min_child_weight=0.001

```

```
min_criterion_penalty=1e5, min_criterion_weight=-1e5),
min_split_gain=0.0, n_estimators=100, n_jobs=-1,
random_state=21, verbose=0, warm_start=False)
```

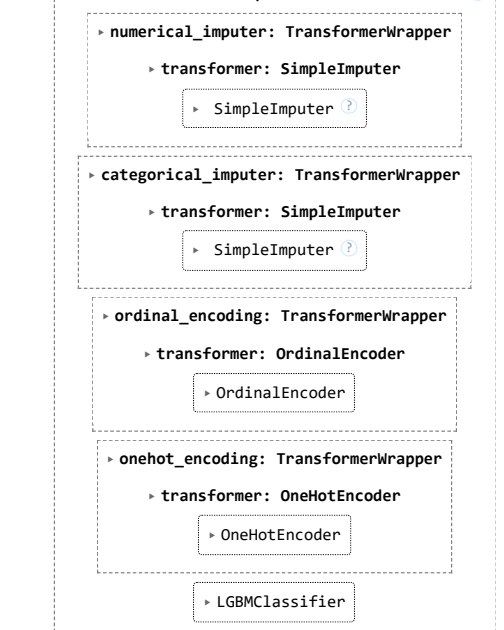
```
reg_alpha=0.0, reg_lambda=0.0, subsample=1.0,  
subsample for bin=200000, subsample freq=8000].
```

```
verbose=false),
'E:\\ICOPD_Model.pd1')
```

```
loaded_best_pipeline = load_model('E:\GPD_Model')
```

loaded\_best\_pipeline

Transformation Pipeline and Model Successfully loaded



## Discussion

The Light Gradient Boosting Machine (LightGBM) model has predicted with precision at 83% accuracy