

Cancel or Not

Import tools

```
1 | import pandas as pd                # basic data tools
2 | import numpy as np
3 | from sklearn.preprocessing import LabelEncoder    # reprocessing tools
4 | from sklearn.preprocessing import StandardScaler
5 | from sklearn.model_selection import GridSearchCV  # model selection
6 | from sklearn import model_selection
7 | from sklearn import metrics
8 | import joblib                          # model saving
9 | import datetime
10 | from sklearn.linear_model import LogisticRegression #Models:# # LogisticRegression
11 | from sklearn.ensemble import RandomForestClassifier    # # Random Forest
12 | import lightgbm as lgbm                            # # LightGBM
13 | from xgboost import XGBClassifier                    # # XGBoost
14 | import seaborn as sns                               # plotting
15 | from matplotlib.pyplot import savefig
16 | from matplotlib.pyplot import axes
17 | from matplotlib import pyplot as plt
```

Accessing data:

```
1 | # df = pd.read_csv('../data/hotel_bookings.csv')
2 | train = pd.read_csv('../data/train.csv') # pre-splitted dataset in previous work
3 | test = pd.read_csv('../data/test.csv')
4 |
5 | # Extra features spreadsheet for hepling understanding
6 | features = pd.read_csv('../data/feature_description.csv', header=None)
7 | features.columns = ['feature', 'description']
```

Cleaning, Exploration and Basic Feature Engineering

Task already done, see `part1.ipynb` if needed.

Preparation

Check features:

```
1 | # show features
2 | pd.set_option("display.max_colwidth", 200)
3 | pd.merge(features, train.describe().T, left_on='feature', right_index=True, how="right")
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

	feature	description	count	mean	std	min	25%	50%	75%
1	is_canceled	order canceled (1) or not (0)	95512.0	0.370414	0.482918	0.00	0.0000	0.0	1.0
2	lead_time	Number of days that elapsed between the entering date of the booking into the PMS and the arrival date	95512.0	104.151583	106.821437	0.00	18.0000	69.0	160.0
5	arrival_date_week_number	Week number of year for arrival date	95512.0	27.171047	13.593740	1.00	16.0000	28.0	38.0
6	arrival_date_day_of_month	Day of arrival date	95512.0	15.797041	8.785365	1.00	8.0000	16.0	23.0
7	stays_in_weekend_nights	Number of weekend nights (Saturday or Sunday) the guest stayed or booked to stay at the hotel	95512.0	0.927852	0.997446	0.00	0.0000	1.0	2.0
8	stays_in_week_nights	Number of week nights (Monday to Friday) the guest stayed or booked to stay at the hotel	95512.0	2.500806	1.905649	0.00	1.0000	2.0	3.0
9	adults	Number of adults	95512.0	1.856866	0.594198	0.00	2.0000	2.0	2.0
10	children	Number of children	95512.0	0.104175	0.399210	0.00	0.0000	0.0	0.0
11	babies	Number of babies	95512.0	0.008020	0.094440	0.00	0.0000	0.0	0.0
16	is_repeated_guest	Value indicating if the booking name was from a repeated guest (1) or not (0)	95512.0	0.031640	0.175041	0.00	0.0000	0.0	0.0
17	previous_cancellations	Number of previous bookings that were cancelled by the customer prior to the current booking	95512.0	0.086617	0.830586	0.00	0.0000	0.0	0.0
18	previous_bookings_not_canceled	Number of previous bookings not cancelled by the customer prior to the current booking	95512.0	0.139899	1.546800	0.00	0.0000	0.0	0.0
21	booking_changes	Number of changes/amendments made to the booking from the moment the booking was entered on the PMS until the moment of check-in or cancellation	95512.0	0.219972	0.649299	0.00	0.0000	0.0	0.0
23	agent	ID of the travel agency that made the booking	82419.0	86.795436	110.798087	1.00	9.0000	14.0	229.0
24	company	ID of the company/entity that made the booking or responsible for paying the booking. ID is presented instead of designation for anonymity reasons	5454.0	189.881371	132.301220	6.00	62.0000	179.0	270.0
25	days_in_waiting_list	Number of days the booking was in the waiting list before it was confirmed to the customer	95512.0	2.307710	17.414238	0.00	0.0000	0.0	0.0

	feature	description	count	mean	std	min	25%	50%	75%
27	adr	Average Daily Rate as defined by dividing the sum of all lodging transactions by the total number of staying nights	95512.0	101.871490	51.137746	-6.38	69.0075	95.0	126.
28	required_car_parking_spaces	Number of car parking spaces required by the customer	95512.0	0.062149	0.245172	0.00	0.0000	0.0	0.0
29	total_of_special_requests	Number of special requests made by the customer (e.g. twin bed or high floor)	95512.0	0.572525	0.793103	0.00	0.0000	0.0	1.0

Use plain integer to encode categorical features according to [official documentation](#) and explicit explanation in sector [Parameters](#).

Preprocessing

Encoding categorical features:

```

1  categoricals = [
2      'hotel',
3      'arrival_date_month',
4      'meal',
5      'country',
6      'market_segment',
7      'distribution_channel',
8      'reserved_room_type',
9      'assigned_room_type',
10     'deposit_type',
11     'agent',
12     'company',
13     'customer_type',
14 ]
15
16 df = train.append(test) # construct a full dataframe for consistent encoding
17 df = df.fillna('0')
18 df['children'] = df['children'].astype(int) # type convertting
19 df_cat_tmp = df[categoricals].astype(str).apply(LabelEncoder().fit_transform) # encode cat columns
20 # df_encoded = pd.merge(df.drop(categoricals, axis=1), df_cat_tmp, left_index=True, right_index=True) # merge numeric and encoded cat columns
21 df[categoricals] = df_cat_tmp[categoricals].astype(int) # merge numeric and encoded cat columns
22 train_encoded = df.iloc[:len(train)]
23 test_encoded = df.iloc[len(train):]

```

Splitting encoded dataset:

```

1  target = 'is_canceled' # splitting features and target
2
3  X_train, X_test = train_encoded.drop(target, axis=1), test_encoded.drop(target, axis=1)
4  y_train, y_test = train_encoded[target], test_encoded[target]

```

Standardization (note that categorical features ignored due to model requirements):

```

1  X_num_train, X_num_test = X_train.drop(categoricals, axis=1), X_test.drop(categoricals, axis=1) # extracting numerical features
2
3  # scaling
4  scaler = StandardScaler()
5  scaler.fit(X_num_train)
6  scaled_features_train = scaler.transform(X_num_train)
7  scaled_features_test = scaler.transform(X_num_test)
8
9  # merge numeric & categorical features
10 scaled_features_df_train = pd.DataFrame(scaled_features_train, index=X_num_train.index, columns=X_num_train.columns)
11 scaled_features_df_test = pd.DataFrame(scaled_features_test, index=X_num_test.index, columns=X_num_test.columns)
12
13 # create standerized dataframes
14 X_train_std = pd.merge(X_train[categoricals], scaled_features_df_train, left_index=True, right_index=True)
15 X_test_std = pd.merge(X_test[categoricals], scaled_features_df_test, left_index=True, right_index=True)

```

Modelling

Define a customized GridSearchCV function for convenience:

```

1  def algorithm_pipeline(model, \

```

```

2         param_grid, \
3         X_train_data = X_train_std, \
4         X_test_data = X_test_std, \
5         y_train_data = y_train, \
6         y_test_data = y_test, \
7         cv=10, \
8         scoring_fit='accuracy',
9         do_probabilities = False):
10
11     gs = GridSearchCV(
12         estimator=model,
13         param_grid=param_grid,
14         cv=cv,
15         n_jobs=-1,
16         scoring=scoring_fit,
17         verbose=2,
18         refit=True # return the best refitted model
19     )
20
21     fitted_model = gs.fit(X_train_data, y_train_data)
22
23     if do_probabilities:
24         y_pred = fitted_model.predict_proba(X_test_data)
25     else:
26         y_pred = fitted_model.predict(X_test_data)
27
28     return fitted_model, y_pred

```

1. Logistic Regrsson (baseline)

```

1 lr = LogisticRegression(penalty='l2', max_iter=7600)
2 lr = lr.fit(X_train, y_train)

```

```

1 y_pred = lr.predict(X_test)
2 print(f'Train(accuracy): {lr.score(X_train, y_train).round(5)}')
3 print(f'Test(accuracy): {lr.score(X_test, y_test).round(5)}\n')
4 print(metrics.classification_report(y_test, y_pred))
5
6 cm_matrix = pd.DataFrame(data=metrics.confusion_matrix(y_test, y_pred), columns=['PP', 'PN'], index=['AP', 'AN'])
7 cm_plot = sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='RdBu_r') # visualize confusion matrix with seaborn heatmap
8 cm_plot

```

```

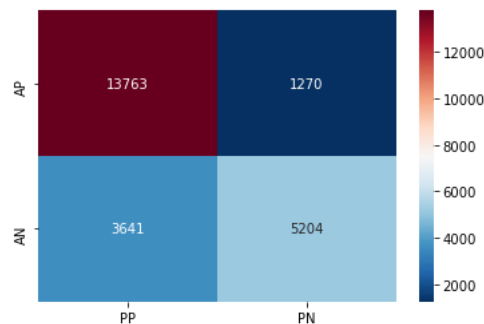
1 Train(accuracy): 0.79377
2 Test(accuracy): 0.79433
3
4      precision    recall  f1-score   support
5
6   0       0.79      0.92      0.85     15033
7   1       0.80      0.59      0.68      8845
8
9  accuracy          0.79      0.79      0.79     23878
10 macro avg       0.80      0.75      0.76     23878
11 weighted avg    0.80      0.79      0.79     23878

```

```

1 <matplotlib.axes._subplots.AxesSubplot at 0x1ac18372448>

```



```

1 clf = lr
2 now = datetime.datetime.now()
3 file_name = "../models/LogisticRegression_" + now.strftime("%m%d%H%M")
4
5 # save parameters
6 with open(f"{file_name}_{gs.best_score_.round(4)*100}%.txt", "w") as para_file:

```

```

7 |     para_file.write(str(gs.best_params_))
8 |     print(gs.best_params_)
9 |
10 | # save figure
11 | cm_plot.get_figure().savefig(file_name + ".png")
12 |
13 | # save model
14 | joblib.dump(clf, file_name + ".model")
15 |
16 | # use this to load a saved model
17 | # lgbm_loaded = joblib.load('.model')

```

```

1 | {'boosting': 'dart', 'feature_fraction': 0.7, 'lambda_l2': 0.1, 'max_depth': 25, 'min_split_gain': 0.1, 'n_estimators': 3000,
  | 'num_leaves': 100, 'objective': 'binary'}

```

```

1 | ['./models/LogisticRegression_04161442.model']

```

2. Random Forest (baseline)

Fit the training set:

```

1 | # feat_labels = X_train.columns
2 | # X_train = X_train[feat_labels]
3 | forest = RandomForestClassifier(n_estimators = 100, max_depth = 25, n_jobs = -1,
4 |                               random_state = 0, bootstrap = True)
5 | forest.fit(X_train, y_train)

```

```

1 | RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
2 |                       criterion='gini', max_depth=25, max_features='auto',
3 |                       max_leaf_nodes=None, max_samples=None,
4 |                       min_impurity_decrease=0.0, min_impurity_split=None,
5 |                       min_samples_leaf=1, min_samples_split=2,
6 |                       min_weight_fraction_leaf=0.0, n_estimators=100,
7 |                       n_jobs=-1, oob_score=False, random_state=0, verbose=0,
8 |                       warm_start=False)

```

Prediction and result:

```

1 | y_pred = forest.predict(X_test)
2 | print(f'Train(accuracy): {forest.score(X_train, y_train).round(5)}')
3 | print(f'Test(accuracy): {forest.score(X_test, y_test).round(5)}\n')
4 | print(metrics.classification_report(y_test, y_pred))
5 |
6 | cm_matrix = pd.DataFrame(data=metrics.confusion_matrix(y_test, y_pred), columns=['PP', 'PN'], index=['AP', 'AN'])
7 | cm_plot = sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='RdBu_r') # visualize confusion matrix with seaborn heatmap
8 | cm_plot

```

```

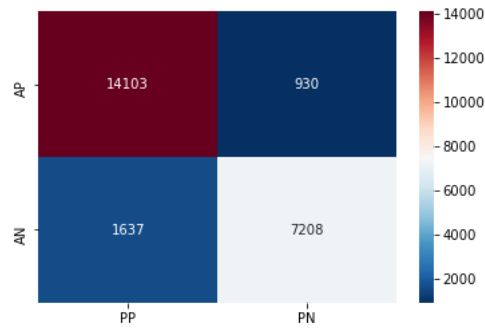
1 | Train(accuracy): 0.97796
2 | Test(accuracy): 0.8925
3 |
4 |           precision    recall  f1-score   support
5 |
6 |    0       0.90       0.94       0.92       15033
7 |    1       0.89       0.81       0.85       8845
8 |
9 |   accuracy                0.89       23878
10 |  macro avg       0.89       0.88       0.88       23878
11 | weighted avg       0.89       0.89       0.89       23878

```

```

1 | <matplotlib.axes._subplots.AxesSubplot at 0x1405cbe3188>

```



Feature importance:

```

1 importances = forest.feature_importances_
2 feat_labels = X_train.columns
3 indices = np.argsort(importances)[::-1]
4 for f in range(X_train.shape[1]):
5     print("%2d) %-*s %f" % (f + 1, 30,
6                             feat_labels[f],
7                             importances[indices[f]]))

```

```

1 1) hotel 0.146114
2 2) lead_time 0.120920
3 3) arrival_date_month 0.115639
4 4) arrival_date_week_number 0.068819
5 5) arrival_date_day_of_month 0.062851
6 6) stays_in_weekend_nights 0.059945
7 7) stays_in_week_nights 0.046318
8 8) adults 0.045873
9 9) children 0.044790
10 10) babies 0.035386
11 11) meal 0.030336
12 12) country 0.026822
13 13) market_segment 0.026390
14 14) distribution_channel 0.025442
15 15) is_repeated_guest 0.024904
16 16) previous_cancellations 0.022886
17 17) previous_bookings_not_canceled 0.019146
18 18) reserved_room_type 0.014447
19 19) assigned_room_type 0.011666
20 20) booking_changes 0.011630
21 21) deposit_type 0.010561
22 22) agent 0.009580
23 23) company 0.005719
24 24) days_in_waiting_list 0.004627
25 25) customer_type 0.004015
26 26) adr 0.002204
27 27) required_car_parking_spaces 0.002110
28 28) total_of_special_requests 0.000859

```

Save:

```

1 clf = forest
2 now = datetime.datetime.now()
3 file_name = "../models/RandomForest_" + now.strftime("%m%d%H%M")
4
5 # save parameters
6 with open(f"{file_name}_{gs.best_score_.round(4)*100}.txt", "w") as para_file:
7     para_file.write(str(gs.best_params_))
8     print(gs.best_params_)
9
10 # save figure
11 cm_plot.get_figure().savefig(file_name + ".png")
12
13 # save model
14 joblib.dump(clf, file_name + ".model")
15
16 # use this to load a saved model
17 # lgbm_loaded = joblib.load('.model')

```

```

1 {'colsample_bytree': 0.7, 'max_depth': 50, 'n_estimators': 100, 'reg_alpha': 1.3, 'reg_lambda': 1.1, 'subsample': 0.9}

```

```

1 ['../models/RandomForest_04160226.model']

```

3. LightGBM

A fast, distributed, high performance gradient boosting (GBT, GBDT, GBRT, GBM or MART) framework based on decision tree algorithms, used for ranking, classification and many other machine learning tasks.

[GitHub](#) | [Documentation](#) | [Parameter Tuning \(official tutorial\)](#)

```
1 # !pip install lightgbm # if not installed
2 lgbmclf = lgbm.LGBMClassifier()
3 param_grid = {
4     # 'objective':['binary','cross_entropy'],
5     'objective':['binary'],
6     'boosting':['gbdt', 'dart'],
7     # 'boosting':['dart'],
8     'n_estimators': [3000],
9     'num_leaves': [25, 50, 100],
10    # 'min_data_in_leaf': []
11    'max_depth': [10, 25],
12    # 'reg_alpha': [1.1, 1.3],
13    # 'reg_lambda': [1.1, 1.3],
14    'lambda_l2': [0.1, 0.3],
15    'min_split_gain': [0.1, 0.3],
16    'feature_fraction': [0.5, 0.7],
17    # 'subsample': [0.7, 0.9],
18    # 'subsample_freq': [20]
19    # 'max_bin': [100, 150]
20 }
21 gs, pred = algorithm_pipeline(lgbmclf, param_grid, cv=3) # some parameters already define above, see function algorithm_pipeline()
22 print(gs.best_score_)
23 print(gs.best_params_)
```

```
1 | Fitting 3 folds for each of 96 candidates, totalling 288 fits
```

```
1 [Parallel(n_jobs=-1)]: Using backend LokyBackend with 6 concurrent workers.
2 [Parallel(n_jobs=-1)]: Done 29 tasks      | elapsed: 1.1min
3 [Parallel(n_jobs=-1)]: Done 150 tasks    | elapsed: 14.9min
4 [Parallel(n_jobs=-1)]: Done 288 out of 288 | elapsed: 205.9min finished
```

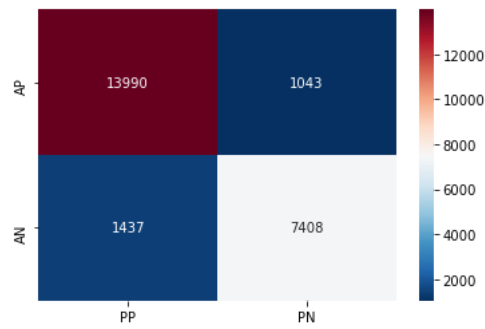
```
1 0.8903593200853003
2 {'boosting': 'dart', 'feature_fraction': 0.7, 'lambda_l2': 0.1, 'max_depth': 25, 'min_split_gain': 0.1, 'n_estimators': 3000,
  'num_leaves': 100, 'objective': 'binary'}
```

Prediction and result:

```
1 clf = gs.best_estimator_
2 clf.fit(X_train, y_train)
3 print(f'Train(accuracy): {clf.score(X_train, y_train).round(5)}') # Python >= 3.7
4 print(f'Test(accuracy): {clf.score(X_test, y_test).round(5)}\n')
5 y_pred = clf.predict(X_test)
6 print(metrics.classification_report(y_test, y_pred))
7
8 cm_matrix = pd.DataFrame(data=metrics.confusion_matrix(y_test, y_pred), columns=['PP', 'PN'], index=['AP', 'AN'])
9 cm_plot = sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='RdBu_r') # visualize confusion matrix with seaborn heatmap
10 cm_plot
```

```
1 Train(accuracy): 0.98896
2 Test(accuracy): 0.89614
3
4          precision    recall  f1-score   support
5
6     0       0.91       0.93       0.92       15033
7     1       0.88       0.84       0.86        8845
8
9    accuracy                0.90       23878
10   macro avg                0.89       0.88       0.89       23878
11   weighted avg                0.90       0.90       0.90       23878
```

```
1 | <matplotlib.axes._subplots.AxesSubplot at 0x1ac5f715f48>
```



Dump the model to file for persistence:

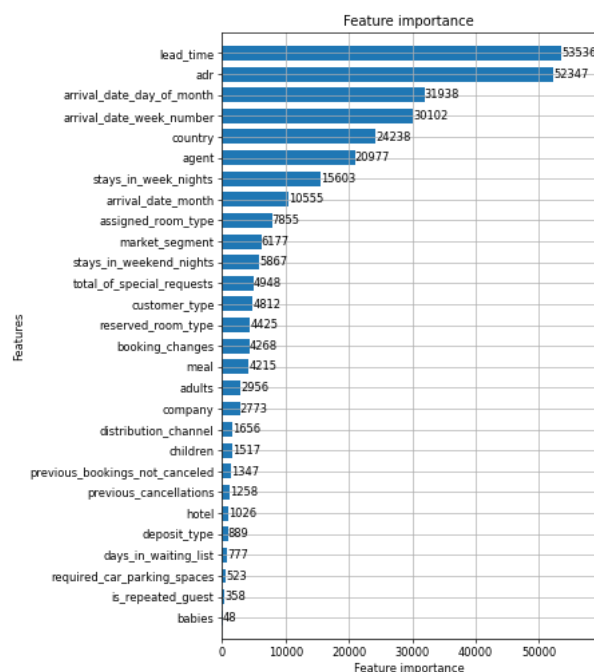
```
1 now = datetime.datetime.now()
2 file_name = "../models/LightGBM_" + now.strftime("%m%d%H%M")
3
4 # save parameters
5 with open(f"{file_name}_{gs.best_score_.round(4)*100}%.txt", "w") as para_file:
6     para_file.write(str(gs.best_params_))
7     print(gs.best_params_)
8
9 # save figure
10 cm_plot.get_figure().savefig(file_name + ".png")
11
12 # save model
13 joblib.dump(clf, file_name + ".model")
14
15 # use this to load a saved model
16 # lgbm_loaded = joblib.load('.model')
```

```
1 {'boosting': 'dart', 'feature_fraction': 0.7, 'lambda_l2': 0.1, 'max_depth': 25, 'min_split_gain': 0.1, 'n_estimators': 3000,
  'num_leaves': 100, 'objective': 'binary'}
```

```
1 ['../models/LightGBM_04161347.model']
```

```
1 lgbm.plot_importance(clf, height = 0.7, figsize = (6, 10), dpi = 60)
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1ac71de8c88>
```




```

1 # !pip install graphviz
2 # lgbm.plot_tree(c1f, figsize = (6, 10), dpi = 60)

```

```

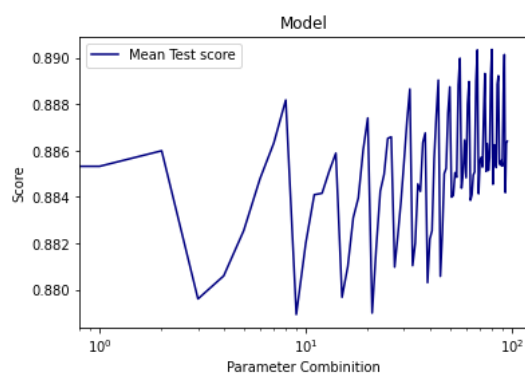
1 Collecting graphviz
2   Downloading graphviz-0.13.2-py2.py3-none-any.whl (17 kB)
3 Installing collected packages: graphviz
4 Successfully installed graphviz-0.13.2

```

```

1 # train_scores_mean = gs.cv_results_["mean_train_score"]
2 # train_scores_std = gs.cv_results_["std_train_score"]
3 test_scores_mean = gs.cv_results_["mean_test_score"]
4 test_scores_std = gs.cv_results_["std_test_score"]
5 params = range(len(gs.cv_results_["params"]))
6
7 plt_cv = plt.figure()
8 plt.title('Model')
9 plt.xlabel('Parameter Combination')
10 plt.ylabel('Score')
11 # plot train scores
12 plt.semilogx(params, test_scores_mean, label='Mean Test score',
13             color='navy')
14 plt.legend(loc='best')
15 plt.show()
16 plt_cv.savefig(file_name + "_cvresult.png")

```



4. XGBoost

XGBoost - Extreme Gradient Boosting, which is an efficient implementation of the gradient boosting framework from Chen & Guestrin (2016) [doi:10.1145/2939672.2939785](https://doi.org/10.1145/2939672.2939785).

[Github](#) | [Documentation](#) | [Parameter Tuning \(official tutorial\)](#)

Parameter tuning using GS

```

1 xgb = XGBClassifier()
2 param_grid = {
3     'n_estimators': [50, 100, 500],
4     'colsample_bytree': [0.7, 0.8],
5     'max_depth': [5, 15, 25, 50],
6     'reg_alpha': [1.1, 1.3],
7     'reg_lambda': [1.1, 1.3],
8     'subsample': [0.7, 0.9],
9     # 'tree_method': ['gpu_hist'],
10    # 'gpu_id': [1]
11 }
12 gs, pred = algorithm_pipeline(xgb, param_grid, cv=3)
13 print(gs.best_score_)
14 print(gs.best_params_)

```

```

1 Fitting 3 folds for each of 192 candidates, totalling 576 fits

```

```

1 [Parallel(n_jobs=-1)]: Using backend LokyBackend with 6 concurrent workers.
2 [Parallel(n_jobs=-1)]: Done 29 tasks      | elapsed: 31.0s
3 [Parallel(n_jobs=-1)]: Done 150 tasks    | elapsed: 16.3min
4 [Parallel(n_jobs=-1)]: Done 353 tasks    | elapsed: 47.1min
5 [Parallel(n_jobs=-1)]: Done 576 out of 576 | elapsed: 91.7min finished

```

```

1 0.8873858846569945
2 {'colsample_bytree': 0.7, 'max_depth': 50, 'n_estimators': 100, 'reg_alpha': 1.3, 'reg_lambda': 1.1, 'subsample': 0.9}

```

Prediction and result:

```

1 clf = gs.best_estimator_
2 clf.fit(X_train, y_train)
3 print(f'Train(accuracy): {clf.score(X_train, y_train).round(5)}') # Python >= 3.7
4 print(f'Test(accuracy): {clf.score(X_test, y_test).round(5)}\n')
5 y_pred = clf.predict(X_test)
6 print(metrics.classification_report(y_test, y_pred))
7
8 cm_matrix = pd.DataFrame(data=metrics.confusion_matrix(y_test, y_pred), columns=['PP', 'PN'], index=['AP', 'AN'])
9 cm_plot = sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='RdBu_r') # visualize confusion matrix with seaborn heatmap
10 cm_plot

```

```

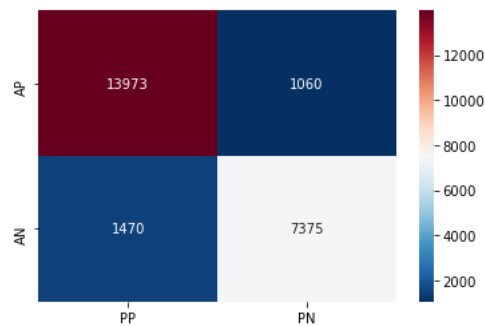
1 Train(accuracy): 0.99574
2 Test(accuracy): 0.89404
3
4      precision    recall  f1-score   support
5
6  0       0.90       0.93       0.92    15033
7  1       0.87       0.83       0.85     8845
8
9  accuracy                0.89    23878
10 macro avg              0.89    0.88    0.89    23878
11 weighted avg          0.89    0.89    0.89    23878

```

```

1 | <matplotlib.axes._subplots.AxesSubplot at 0x1405e6b0248>

```



Dumping to file:

```

1 now = datetime.datetime.now()
2 file_name = "../models/XGBoost_" + now.strftime("%m%d%H%M")
3
4 # save parameters
5 with open(f"{file_name}_{gs.best_score_.round(4)*100}%.txt", "w") as para_file:
6     para_file.write(str(gs.best_params_))
7     print(gs.best_params_)
8
9 # save figure
10 cm_plot.get_figure().savefig(file_name + ".png")
11
12 # save model
13 joblib.dump(clf, file_name + ".model")
14
15 # use this to load a saved model
16 # lgbm_loaded = joblib.load('.model')

```

```

1 | {'colsample_bytree': 0.7, 'max_depth': 50, 'n_estimators': 100, 'reg_alpha': 1.3, 'reg_lambda': 1.1, 'subsample': 0.9}

```

```

1 | ['../models/XGBoost_04160226.model']

```