

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 # Loading datasets
        4 train_df = pd.read_csv("training.csv")
        5 test_df = pd.read_csv("test.csv")
        6 agree_df = pd.read_csv("check_agreement.csv")
        7 corr_df = pd.read_csv("check_correlation.csv")
```

writing functions for agreement and correlation test

```

In [2]: 1 # Check agreement test
2
3
4 # agreement test as mentioned in kaggle resources#
5 from sklearn.metrics import roc_curve, auc
6
7 def __roc_curve_split(data_zero, data_one, sample_weights_zero, sample_weights_one):
8     """
9     Compute roc curve
10
11     :param data_zero: 0-labeled data
12     :param data_one: 1-labeled data
13     :param sample_weights_zero: weights for 0-labeled data
14     :param sample_weights_one: weights for 1-labeled data
15     :return: roc curve
16     """
17     labels = [0] * len(data_zero) + [1] * len(data_one)
18     weights = np.concatenate([sample_weights_zero, sample_weights_one])
19     data_all = np.concatenate([data_zero, data_one])
20     fpr, tpr, _ = roc_curve(labels, data_all, sample_weight=weights)
21     return fpr, tpr
22
23 def compute_ks(data_prediction, mc_prediction, weights_data, weights_mc):
24     """
25     Compute Kolmogorov-Smirnov (ks) distance between real data predictions and Monte Carlo predictions
26
27     :param data_prediction: array-like, real data predictions
28     :param mc_prediction: array-like, Monte Carlo data predictions
29     :param weights_data: array-like, real data weights
30     :param weights_mc: array-like, Monte Carlo weights
31     :return: ks value
32     """
33     assert len(data_prediction) == len(weights_data), 'Data length and weight length must be equal'
34     assert len(mc_prediction) == len(weights_mc), 'Data length and weight length must be equal'
35
36     data_prediction, mc_prediction = np.array(data_prediction), np.array(mc_prediction)
37     weights_data, weights_mc = np.array(weights_data), np.array(weights_mc)
38
39     assert np.all(data_prediction >= 0.) and np.all(data_prediction <= 1.), 'Data predictions must be between 0 and 1'
40     assert np.all(mc_prediction >= 0.) and np.all(mc_prediction <= 1.), 'MC predictions must be between 0 and 1'
41
42     weights_data /= np.sum(weights_data)
43     weights_mc /= np.sum(weights_mc)
44
45     fpr, tpr = __roc_curve_split(data_prediction, mc_prediction, weights_data, weights_mc)
46
47     Dnm = np.max(np.abs(fpr - tpr))
48     return Dnm
49
50 # check correlation test
51
52 # correlation test as mentioned in kaggle resources
53
54 def __rolling_window(data, window_size):
55     """
56     Rolling window: take window with definite size through the array

```

```

57
58     :param data: array-like
59     :param window_size: size
60     :return: the sequence of windows
61
62     Example: data = array(1, 2, 3, 4, 5, 6), window_size = 4
63     Then this function return array(array(1, 2, 3, 4), array(2, 3, 4, 5))
64     """
65     shape = data.shape[:-1] + (data.shape[-1] - window_size + 1, window_size)
66     strides = data.strides + (data.strides[-1],)
67     return np.lib.stride_tricks.as_strided(data, shape=shape, strides=strides)
68
69 def __cvm(subindices, total_events):
70     """
71     Compute Cramer-von Mises metric.
72     Compared two distributions, where first is subset of second one.
73     Assuming that second is ordered by ascending
74
75     :param subindices: indices of events which will be associated with the first distribution
76     :param total_events: count of events in the second distribution
77     :return: cvm metric
78     """
79     target_distribution = np.arange(1, total_events + 1, dtype='float') / total_events
80     subarray_distribution = np.cumsum(np.bincount(subindices, minlength=total_events)) / total_events
81     subarray_distribution /= 1.0 * subarray_distribution[-1]
82     return np.mean((target_distribution - subarray_distribution) ** 2)
83
84 def compute_cvm(predictions, masses, n_neighbours=200, step=50):
85     """
86     Computing Cramer-von Mises (cvm) metric on background events: take average of cvm metric for each mass bin
87     In each mass bin global prediction's cdf is compared to prediction's cdf
88
89     :param predictions: array-like, predictions
90     :param masses: array-like, in case of Kaggle tau23mu this is reconstructed mass
91     :param n_neighbours: count of neighbours for event to define mass bin
92     :param step: step through sorted mass-array to define next center of bin
93     :return: average cvm value
94     """
95     predictions = np.array(predictions)
96     masses = np.array(masses)
97     assert len(predictions) == len(masses)
98
99     # First, reorder by masses
100    predictions = predictions[np.argsort(masses)]
101
102    # Second, replace probabilities with order of probability among other events
103    predictions = np.argsort(np.argsort(predictions, kind='mergesort'), kind='mergesort')
104
105    # Now, each window forms a group, and we can compute contribution of each window
106    cvms = []
107    for window in __rolling_window(predictions, window_size=n_neighbours):
108        cvms.append(__cvm(subindices=window, total_events=len(predictions)))
109    return np.mean(cvms)

```

function to add new features

```

In [3]: 1 # feature engineering
2 def new_feats(df):
3     df2 = df.copy()
4     df2['isolation_abc'] = df['isolationa'] + df['isolationb'] + df['isolationc']
5     df2['isolation_def'] = df['isolationd'] + df['isolatione'] + df['isolationf']
6     df2['p_IP'] = df['p0_IP'] + df['p1_IP'] + df['p2_IP']
7     df2['p_p'] = df['p0_p'] + df['p1_p'] + df['p2_p']
8     df2['IP_pp'] = df['IP_p0p2'] + df['IP_p1p2']
9     df2['p_IPSig'] = df['p0_IPSig'] + df['p1_IPSig'] + df['p2_IPSig']
10    #new feature using 'FlightDistance' and LifeTime(from literature)
11    df2['FD_LT'] = df['FlightDistance'] / df['LifeTime']
12    #new feature using 'FlightDistance', 'p0_p', 'p1_p', 'p2_p'(from literature)
13    df2['FD_p0p1p2_p'] = df['FlightDistance'] / (df['p0_p'] + df['p1_p'] + df['p2_p'])
14    #new feature using 'LifeTime', 'p0_IP', 'p1_IP', 'p2_IP'(from literature)
15    df2['NEW5_lt'] = df['LifeTime'] * (df['p0_IP'] + df['p1_IP'] + df['p2_IP']) / 3
16    #new feature using 'p0_track_Chi2Dof', 'p1_track_Chi2Dof', 'p2_track_Chi2Dof'
17    df2['Chi2Dof_MAX'] = df.loc[:, ['p0_track_Chi2Dof', 'p1_track_Chi2Dof', 'p2_track_Chi2Dof']].max(axis=1)
18    # features from kaggle discussion forum
19    df2['flight_dist_sig2'] = (df['FlightDistance'] / df['FlightDistanceError']) ** 2
20    df2['flight_dist_sig'] = df['FlightDistance'] / df['FlightDistanceError']
21    df2['NEW_IP_dira'] = df['IP'] * df['dira']
22    df2['p0p2_ip_ratio'] = df['IP'] / df['IP_p0p2']
23    df2['p1p2_ip_ratio'] = df['IP'] / df['IP_p1p2']
24    df2['DCA_MAX'] = df.loc[:, ['DOCAone', 'DOCAtwo', 'DOCAthree']].max(axis=1)
25    df2['iso_bdt_min'] = df.loc[:, ['p0_IsoBDT', 'p1_IsoBDT', 'p2_IsoBDT']].min(axis=1)
26    df2['iso_min'] = df.loc[:, ['isolationa', 'isolationb', 'isolationc', 'isolationd', 'isolatione', 'isolationf']].min(axis=1)
27    return df2
28

```

```

In [4]: 1 # adding engineered features to training and test datasets
2 train_df_1 = new_feats(train_df)
3 test_df_1 = new_feats(test_df)

```

```

In [5]: 1 # identifying some features to remove which have been used to engineer new features
2 #Low importance in EDA
3 remove = ['id', 'min_ANNmuon', 'production', 'mass', 'signal', 'SPDhits', 'CDF', 'p0_pt', 'p1_pt', 'p2_pt', 'p0_p', 'p1_p', 'p2_p', 'p0_eta', 'p1_eta', 'p2_eta', 'isolationc', 'isolationd', 'isolatione', 'isolationf', 'p0_IsoBDT', 'p1_IsoBDT', 'p2_IsoBDT', 'IP_p0p2', 'IP_p1p2', 'p0_track_Chi2Dof', 'p1_track_Chi2Dof', 'p2_track_Chi2Dof', 'p0_IPSig', 'p1_IPSig', 'p2_IPSig', 'DOCAone', 'DOCAtwo', 'DOCAthree']
4
5
6
7
8 # making a list of features to be used to train the model and make predictions
9 features = list(f for f in train_df_1.columns if f not in remove)

```

```

In [6]: 1 len(features)

```

Out[6]: 28

Creating a new class for UGradientBoosting with loss incorporated in the class itself for bayesian optimization hyper-parameter tuning

```

In [7]: 1 from hep_ml.gradientboosting import UGradientBoostingClassifier
2 from hep_ml.losses import BinFlatnessLossFunction
3 from sklearn.base import BaseEstimator
4 from sklearn.metrics import accuracy_score
5 from collections import Counter
6 class UGradientBoostingClassifierWithLoss(BaseEstimator):
7     def __init__(
8         # self, max_depth=3, max_features=0.8, learning_rate=0.01,
9         # n_estimators=80, subsample=0.8
10        self, max_depth, n_estimators, **params
11    ):
12        loss = BinFlatnessLossFunction(
13            ['mass'], n_bins=15, uniform_label = 0 , fl_coefficient=15, power
14        )
15
16        self.estimator = UGradientBoostingClassifier(
17            loss=loss,
18            train_features = list(f for f in train_df_1.columns if f not in
19            max_depth = max_depth,
20            n_estimators = n_estimators,
21            **params
22        )
23
24    def fit(self, X, y=None):
25        self.estimator.fit(X, y)
26
27        return self
28
29    def predict_proba(self, X):
30        return self.estimator.predict_proba(X)
31
32    def predict(self, X):
33        return self.estimator.predict(X)
34
35    def transform(self, X):
36        return self.estimator.transform(X)
37
38    def get_params(self, deep=True):
39        # suppose this estimator has parameters "alpha" and "recursive"
40        params_to_keep = [
41            "max_depth",
42            "max_features",
43            "learning_rate",
44            "n_estimators",
45            "subsample",
46        ]
47
48        ret = dict()
49        tret = self.estimator.get_params(deep=deep)
50        for key in params_to_keep:
51            ret[key] = tret[key]
52
53        return ret
54
55    def set_params(self, **parameters):
56        self.estimator.set_params(parameters)

```

```
57
58     return self
59
60     def score(self, X, y):
61         y_pred = self.estimator.predict(X)
62
63         acc = accuracy_score(y, y_pred)
64         print(acc)
65         print(Counter(y))
66         print(Counter(y_pred))
67
68     return acc
69
```

In [8]:

```
1 from sklearn.model_selection import cross_val_score
2 from sklearn.model_selection import StratifiedKFold
3
4 from hep_ml.gradientboosting import UGradientBoostingClassifier
5 from hep_ml.losses import BinFlatnessLossFunction
6 from bayes_opt import BayesianOptimization
7 from sklearn.metrics import roc_auc_score, make_scorer, accuracy_score
8 import time
9 import json
10 auc = make_scorer(roc_auc_score)
```

bayesian optimization demands bounds of hyperparameters. Therefore parameters with discrete values can not be tuned with this technique. The discrete hyperparameters need to be tuned are `n_estimators` and `max_depth`. Therefore creating all possible pairs of these two parameters and doing grid search for these two parameters while running bayesian optimization for each pair of these two discrete parameters

In [23]:

```

1 start = time.time()
2 # we will save best parameters in best dictionary
3 best = dict()
4 best['score'] = 0
5 import itertools
6 max_depth = [3,6,9]
7 n_estimators = [50,400,900]
8 for x in itertools.product(max_depth, n_estimators):
9
10     print("max_depth-n_estimator pair: {}".format(x))
11     def gbm_cl_bo(max_features, learning_rate, subsample):
12
13         params_gbm = {}
14         # params_gbm['max_depth'] = round(max_depth)
15         params_gbm['max_features'] = max_features
16         params_gbm['learning_rate'] = learning_rate
17         # params_gbm['n_estimators'] = round(n_estimators)
18         params_gbm['subsample'] = subsample
19
20
21         # Loss = BinFlatnessLossFunction(['mass'], n_bins=15, uniform_label
22         skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
23         skf.get_n_splits(train_df_1[features + ['mass']], train_df_1['signal'])
24
25         scores = cross_val_score(UGradientBoostingClassifierWithLoss(max_depth,
26                               train_df_1[features + ['mass']], train_df_1['signal'],
27                               scoring = auc,
28                               cv=skf))
29
30         score = scores.mean()
31         #print("max_depth: {}, n_estimators: {}, max_features: {}, learning_rate: {}, subsample: {}".format(x[0], x[1], params_gbm['max_features'], params_gbm['learning_rate'], params_gbm['subsample']))
32
33         if score > best['score']:
34             best['score'] = score
35             best['max_depth'] = x[0]
36             best['n_estimators'] = x[1]
37             best['max_features'] = params_gbm['max_features']
38             best['learning_rate'] = params_gbm['learning_rate']
39             best['subsample'] = params_gbm['subsample']
40
41         return score
42
43     params_gbm = {
44         # 'max_depth': (3),
45         'max_features': (0.5, 1),
46         'learning_rate': (0.01, 1),
47         # 'n_estimators': (100),
48         'subsample': (0.5, 1)
49     }
50
51     gbm_bo = BayesianOptimization(gbm_cl_bo, params_gbm, random_state=111)
52     gbm_bo.maximize(init_points=10, n_iter=3)
53
54
55
56

```

```
57 print('It takes %s minutes' % ((time.time() - start)/60))
```

```
=====
max_depth-n_estimator pair: (9, 900)
|  iter   |  target  | learni... | max_fe... | subsample |
|-----|-----|-----|-----|-----|
|  1      |  0.8736  |  0.616    |  0.5845   |  0.718    |
|  2      |  0.8736  |  0.7716   |  0.6477   |  0.5746   |
|  3      |  0.8706  |  0.03225  |  0.7101   |  0.6193   |
|  4      |  0.8811  |  0.3443   |  0.9954   |  0.6189   |
|  5      |  0.876   |  0.09038  |  0.8348   |  0.8106   |
|  6      |  0.8793  |  0.2815   |  0.7331   |  0.5592   |
|  7      |  0.8758  |  0.08322  |  0.9504   |  0.897    |
|  8      |  0.8681  |  0.8422   |  0.9076   |  0.9955   |
|  9      |  0.8762  |  0.5815   |  0.9069   |  0.7107   |
| 10      |  0.8722  |  0.03717  |  0.7271   |  0.5527   |
| 11      |  0.8818  |  0.4855   |  0.9904   |  0.504    |
| 12      |  0.877   |  0.0682   |  0.999    |  0.5176   |
| 13      |  0.8817  |  0.4756   |  0.9829   |  0.5015   |
=====
It takes 1945.4509294231732 minutes
```

took nearly 35 hours to run!!

```
max_depth-n_estimator pair: (6, 900)
```

```
iter | target | learni... | max_fe... | subsample |
```

```
12 | 0.8887 | 0.4592 | 0.9965 | 0.5016 |
```

```
4 | 0.8866 | 0.3443 | 0.9954 | 0.6189
```

```
13 | 0.8867 | 0.3242 | 0.9863 | 0.5006
```

```
2 | 0.881 | 0.7716 | 0.6477 | 0.5746
```

```
In [24]: 1 best
```

```
Out[24]: {'score': 0.8886783798768416,
          'max_depth': 6,
          'n_estimators': 900,
          'max_features': 0.9965346255608354,
          'learning_rate': 0.459200553861577,
          'subsample': 0.501638892280019}
```

it is found that the hyperparameter which gives best score does not pass correlation test. Therefore trying other hyperparameters with best scores in descending order

```
In [16]: 1 UGradientBoostingClassifier?
```

following model with passed parameters passes both the required tests. Therefore using the results from this model to check on kaggle test data


```
In [17]: 1 loss = BinFlatnessLossFunction(['mass'], n_bins=15, uniform_label=0 , fl_coe
2 model = UGradientBoostingClassifier(loss=loss, n_estimators=900,
3                                     max_depth = 6,
4                                     learning_rate = 0.1,
5                                     train_features = features,
6                                     subsample=0.5)
7 model.fit(train_df_1[features + ['mass']], train_df_1['signal'])
```

```
Out[17]: UGradientBoostingClassifier(learning_rate=0.1,
                                     loss=BinFlatnessLossFunction(allow_wrong_signs=True,
                                     e,
                                     fl_coefficient=15,
                                     n_bins=15, power=2,
                                     uniform_features=['mass'],
                                     uniform_label=array
                                     ([0])),
                                     max_depth=6, max_features=None, max_leaf_nodes=None,
                                     e,
                                     min_samples_leaf=1, min_samples_split=2,
                                     n_estimators=900,
                                     random_state=RandomState(MT19937) at 0x290F6367740,
                                     splitter=...
                                     train_features=['LifeTime', 'dira',
                                     'FlightDistance',
                                     'FlightDistanceError', 'IP',
                                     'IPSig', 'VertexChi2', 'pt', 'iso',
                                     'ISO_SumBDT', 'isolation_abc',
                                     'isolation_def', 'p_IP', 'p_p',
                                     'IP_pp', 'p_IPSig', 'FD_LT',
                                     'FD_p0p1p2_p', 'NEW5_lt',
                                     'Chi2Dof_MAX', 'flight_dist_sig2',
                                     'flight_dist_sig', 'NEW_IP_dira',
                                     'p0p2_ip_ratio', 'p1p2_ip_ratio',
                                     'DCA_MAX', 'iso_bdt_min',
                                     'iso_min'],
                                     update_tree=True)
```

```
In [22]: 1 # saving the model to the memory
2 import pickle
3 filename = 'finalized_model_1.sav'
4 pickle.dump(model, open(filename, 'wb'))
```

```
In [23]: 1 # Loading the model from memory
2 loaded_model = pickle.load(open(filename, 'rb'))
```

In [24]: 1 loaded_model

```
Out[24]: UGradientBoostingClassifier(learning_rate=0.1,
                                     loss=BinFlatnessLossFunction(allow_wrong_signs=True,
                                                                    fl_coefficient=15,
                                                                    n_bins=15, power=2,
                                                                    uniform_features=['mass'],
                                                                    uniform_label=array
                                                                    ([0])),
                                     max_depth=6, max_features=None, max_leaf_nodes=None,
                                     min_samples_leaf=1, min_samples_split=2,
                                     n_estimators=900,
                                     random_state=RandomState(MT19937) at 0x290A3B58140,
                                     splitter=...,
                                     train_features=['LifeTime', 'dira',
                                                       'FlightDistance',
                                                       'FlightDistanceError', 'IP',
                                                       'IPSig', 'VertexChi2', 'pt', 'iso',
                                                       'ISO_SumBDT', 'isolation_abc',
                                                       'isolation_def', 'p_IP', 'p_p',
                                                       'IP_pp', 'p_IPSig', 'FD_LT',
                                                       'FD_p0p1p2_p', 'NEW5_lt',
                                                       'Chi2Dof_MAX', 'flight_dist_sig2',
                                                       'flight_dist_sig', 'NEW_IP_dira',
                                                       'p0p2_ip_ratio', 'p1p2_ip_ratio',
                                                       'DCA_MAX', 'iso_bdt_min',
                                                       'iso_min'],
                                     update_tree=True)
```

```
In [19]: 1 # conducting agreement check test
2 check_agreement = pd.read_csv("check_agreement.csv")
3 check_agreement = new_feats(check_agreement)
4
5 #check_agreement = pandas.read_csv(folder + 'check_agreement.csv', index_col=0)
6 agreement_probs = model.predict_proba(check_agreement[features])[:, 1]
7
8 ks = compute_ks(
9     agreement_probs[check_agreement['signal'].values == 0],
10    agreement_probs[check_agreement['signal'].values == 1],
11    check_agreement[check_agreement['signal'] == 0]['weight'].values,
12    check_agreement[check_agreement['signal'] == 1]['weight'].values)
13 #print 'KS metric', ks, ks < 0.09
14 print("KS metric {}".format(ks))
15 print(ks < 0.09)
```

KS metric 0.0795798778412874

True

```
In [18]: 1 # conducting the correlation test
2 #check_correlation = pandas.read_csv(folder + 'check_correlation.csv', index_
3 check_correlation = pd.read_csv("check_correlation.csv", index_col = "id")
4 check_correlation = new_feats(check_correlation)
5 correlation_probs = model.predict_proba(check_correlation[features])[:, 1]
6 cvm = compute_cvm(correlation_probs, check_correlation['mass'])
7 #print 'CvM metric', cvm, cvm < 0.002
8 print("CvM metric {}".format(cvm))
9 print(cvm < 0.002)
```

CvM metric 0.00135303046269074

True

```
In [21]: 1 # making submission file with test dataset to be submitted on kaggle
2 test_probs = model.predict_proba(test_df_1[features])[:,1]
3 result = pd.DataFrame({"id": test_df["id"], "prediction": test_probs})
4 result.to_csv("final_result_1.csv", index=False)
```

20th rank with final_result_1.csv

```
In [ ]: 1
```