

In [1]:

```
1 #!/usr/bin/env python
2  # -*- coding: utf-8 -*-
3
4  import os
5  import sys
6  import json
7  import logging
8
9  sys.path.append('/home/moegwjawel/atec_project/train/lmf/lib/python3.6/site-pack
10
11  import numpy as np
12  import pandas as pd
13
14  from sklearn.preprocessing import StandardScaler
15  from sklearn.model_selection import train_test_split
16  from sklearn.model_selection import GridSearchCV
17
18  import joblib
19
20  input_data_path = '/home/admin/workspace/job/input/train.jsonl'
21  output_model_path = '/home/admin/workspace/job/output/train.model'
22  result_path = '/home/admin/workspace/job/output/result.json'
23
24  input_data_path = '/mnt/atec/train.jsonl'
25
26  logging.warning("Begin Train.")
27  # 1. 读取训练数据进行训练
28  import time
29  time_start = time.time()
30  # pd_raw_data = pd.read_json(input_data_path, encoding='utf-8', lines=True, nro
31  pd_raw_data = pd.read_json(input_data_path, encoding='utf-8', lines=True)
32  time_end=time.time()
33  print('time cost', time_end - time_start, 's')
34  logging.warning("Finish Read.")
```

WARNING:root:Begin Train.

WARNING:root:Finish Read.

time cost 20.96392798423767 s

In [2]:

```
1 # 2. 数据预处理
2 pd_data = pd_raw_data.copy(deep=True)
3 print(pd_data.shape)
4 # pd_data = pd_data.drop(pd_data.columns[np.where((pd_data.dtypes!='int64').valu
5 # 2.1 去除 label = -1
6 pd_data = pd_data.drop(pd_data[(pd_data['label'] == -1)].index)
7 # 2.2 去除文字列
8 pd_data = pd_data.drop(columns=['memo_polish'],axis=1)
9 print(pd_data.shape)
10 # 2.3 去除离群点
11 min_max_map={}
12 columns = pd_data.columns.values.tolist()[1:-1] # 去除 id 和 label 的所有列
13 for col in columns:
14     if pd_data.dtypes[col]=='float64':
15         d_min, d_max = pd_data[col].quantile([0.0001, 0.9999])
16     elif pd_data.dtypes[col]=='int64':
17         d_min, d_max = pd_data[col].quantile([0.001, 0.999])
18     else:
19         print("ERROR: =====")
20     print(col,d_min,d_max)
21     min_max_map[col]=(d_min,d_max)
22     pd_data = pd_data.drop(pd_data[(pd_data[col]>d_max)].index)
23     pd_data = pd_data.drop(pd_data[(pd_data[col]<d_min)].index)
24 #     pd_data[pd_data[col]>d_max]=np.nan
25 #     pd_data[pd_data[col]<d_min]=np.nan
26 print(pd_data.shape)
27 # 2.4 去除取值只有一个值的列
28 columns = pd_data.columns.values.tolist()[1:-1] # 去除 id 和 label 的所有列
29 drop_list = []
30 for col in columns:
31     num = len(pd_data[col].unique())
32     if num == 1:
33         drop_list.append(col)
34 pd_data = pd_data.drop(columns=drop_list)
35 print(pd_data.shape)
x353 -1.001 1.0
x354 0.0 116.532000000000652
x355 1.0 51.31267096773692
x356 -100.0 91.0
x357 -1.001 1.0
x358 -100.0 1082.7211013887252
x359 -100.0 392.0
x360 -1.001 1.0
x361 -100.0 1.0
x362 -1.003 0.0
x363 -100.0 1.0
x364 -100.0 50.0
x365 -100.0 0.0
x366 -100.0 36.0
x367 -100.0 -100.0
x368 -1.0 1.0
x369 0.0 161.0
x370 -1.001 1.0
x371 -100.0 -100.0
x372 0.0 1.0
```

In [4]:

```
1 print(len(drop_list))
2 print(drop_list)
3 print(min_max_map)
4 file = open('min_max_map.txt', 'w')
5 for k,v in min_max_map.items():
6     file.write(str(k)+' '+str(v[0]) + ' ' +str(v[1])+'\n')
7 file.close()
```

```
82
['x2', 'x8', 'x13', 'x18', 'x25', 'x32', 'x33', 'x42', 'x48', 'x55',
'x68', 'x70', 'x71', 'x72', 'x77', 'x79', 'x91', 'x96', 'x97', 'x100',
'x107', 'x110', 'x111', 'x113', 'x121', 'x130', 'x135', 'x141', 'x14
5', 'x152', 'x155', 'x178', 'x184', 'x198', 'x202', 'x207', 'x209', 'x
215', 'x224', 'x225', 'x228', 'x232', 'x239', 'x248', 'x249', 'x258',
'x261', 'x264', 'x265', 'x277', 'x279', 'x280', 'x281', 'x289', 'x29
2', 'x300', 'x319', 'x325', 'x341', 'x343', 'x347', 'x351', 'x367', 'x
371', 'x373', 'x384', 'x385', 'x389', 'x394', 'x400', 'x419', 'x427',
'x431', 'x436', 'x441', 'x442', 'x451', 'x452', 'x456', 'x461', 'x46
2', 'x475']
{'x0': (0.0, 1.0), 'x1': (-2.0, 179.0), 'x2': (-100.0, -100.0), 'x3':
(-1.001, 2623.9703999996324), 'x4': (-1.001, 300.0), 'x5': (-1.001, 30
0.0), 'x6': (-100.0, 0.0), 'x7': (-100.0, 100000.0), 'x8': (-1.0, 1.
0), 'x9': (-100.0, 1.0), 'x10': (0.0, 23.0), 'x11': (-1.001, 2.4516447
24085354), 'x12': (0.0, 136.152000000003097), 'x13': (0.0, 2.0), 'x14':
(2.0, 8.0), 'x15': (-1.0, 112.634000000000547), 'x16': (-1.0, 25.0), 'x
17': (0.0, 2030.72500000000058), 'x18': (-1.0, 5113.980449986075), 'x1
9': (-100.0, 465.0), 'x20': (-1.0, 1526.0), 'x21': (-100.0, 1.0), 'x2
```

In [16]:

```
1 # 3. 数据转 numpy
2 np_data = pd_data.values
3 # 去除索引列
4 np_data = np_data[:, 1:]
5 # 分成 input 和 output, output 降维到一维
6 np_data_input = np_data[:, :np_data.shape[1]-1]
7 np_data_input = np_data_input.astype(np.float64)
8
9 print(np_data_input.shape)
10 print(np_data_input)
11
12 np_data_output = np_data[:, np_data.shape[1]-1:]
13 np_data_output = np_data_output.astype(np.int64)
14 np_data_output = np.squeeze(np_data_output, axis=1)
15
16 print(np_data_output.shape)
17 print(np_data_output)
18
19 test_size=0.2
20 X_train, X_test, Y_train, Y_test = train_test_split(np_data_input, np_data_output,
21 logging.warning("Finish Process.")
```

```
(43031, 398)
[[ 1.      -2.      0.      ... -100.      0.      1.    ]
 [ 1.      90.      0.      ... -100.      0.      1.    ]
 [ 1.      69.      0.      ... -100.      0.      7.    ]
 ...
 [ 1.      -2.      2.      ... -100.      0.      1.    ]
 [ 1.     179.      0.      ... -100.     60.42 -100.   ]
 [ 1.      -2.     -1.001 ... -100.      0.      7.    ]]
(43031,)
[0 1 1 ... 0 1 0]
```

WARNING:root:Finish Process.

In [17]:

```
1 # # 4. 数据归一化
2 # sc = StandardScaler()
3 # sc.fit(X_train)
4 # X_train = sc.transform(X_train)
5 # X_test = sc.transform(X_test)
```

In [18]:

```
1 #####
2 logging.warning("Begin Model.")
3 result = {} # 保存一些结果
4 # lightGBM 模型
5 model_name = 'lightGBM'
6 version = 'v0.4_test'
7 from lightgbm import LGBMClassifier
8
9 # 5.1 模型定义
10 model = LGBMClassifier(
11     verbose=1
12 )
13
14 # 5.2 模型训练
15 import time
16 time_start=time.time()
17 model.fit(X_train, Y_train)
18 time_end=time.time()
19 print('time cost',time_end-time_start,'s')
20 # 5.3 模型评估
21 result['score_train'] = model.score(X_train, Y_train)
22 result['score_test'] = model.score(X_test, Y_test)
23 # 5.4 模型保存
24 joblib.dump(model, f'./model/%s_%s_[%f, %f, %f, %s].model' %
25             (model_name, version,
26              result['score_test'], result['score_train'],
27              test_size, model))
28 # 5.5 结果保存
29 print(result)
30 #####
```

WARNING:root:Begin Model.

```
[LightGBM] [Info] Number of positive: 12297, number of negative: 22127
[LightGBM] [Warning] Auto-choosing row-wise multi-threading, the overh
ead of testing was 0.135918 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 35246
[LightGBM] [Info] Number of data points in the train set: 34424, numbe
r of used features: 395
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.357222 -> initscore=
-0.587443
[LightGBM] [Info] Start training from score -0.587443
time cost 3.3233299255371094 s
{'score_train': 0.8687834069254009, 'score_test': 0.8416405251539445}
```

In [19]:

```
1  # #####
2  # logging.warning("Begin Model.")
3  # result = {} # 保存一些结果
4  # # CatBoost 模型
5  # model_name = 'CatBoost'
6  # version = 'v0.4_test'
7  # from catboost import CatBoostClassifier
8
9  # # 5.1 模型定义
10 # model = CatBoostClassifier(
11 #     verbose=1
12 # )
13
14 # # 5.2 模型训练
15 # import time
16 # time_start=time.time()
17 # model.fit(X_train, Y_train)
18 # time_end=time.time()
19 # print('time cost',time_end-time_start,'s')
20 # # 5.3 模型评估
21 # result['score_train'] = model.score(X_train, Y_train)
22 # result['score_test'] = model.score(X_test, Y_test)
23 # # 5.4 模型保存
24 # joblib.dump(model, f'./model/%s_%s_[%f, %f, %f, %s].model' %
25 #             (model_name, version,
26 #              result['score_test'], result['score_train'],
27 #              test_size, model))
28 # # 5.5 结果保存
29 # print(result)
30 # #####
```

```
848:      learn: 0.2919997      total: 19s      remaining: 3.38s
849:      learn: 0.2919033      total: 19s      remaining: 3.36s
850:      learn: 0.2918047      total: 19.1s    remaining: 3.33s
851:      learn: 0.2916512      total: 19.1s    remaining: 3.31s
852:      learn: 0.2915599      total: 19.1s    remaining: 3.29s
853:      learn: 0.2914986      total: 19.1s    remaining: 3.27s
854:      learn: 0.2914383      total: 19.1s    remaining: 3.25s
855:      learn: 0.2913382      total: 19.2s    remaining: 3.22s
856:      learn: 0.2912935      total: 19.2s    remaining: 3.2s
857:      learn: 0.2912220      total: 19.2s    remaining: 3.18s
858:      learn: 0.2911435      total: 19.2s    remaining: 3.15s
859:      learn: 0.2910624      total: 19.2s    remaining: 3.13s
860:      learn: 0.2909642      total: 19.3s    remaining: 3.11s
861:      learn: 0.2908532      total: 19.3s    remaining: 3.09s
862:      learn: 0.2907623      total: 19.3s    remaining: 3.07s
863:      learn: 0.2906677      total: 19.3s    remaining: 3.04s
864:      learn: 0.2905833      total: 19.4s    remaining: 3.02s
865:      learn: 0.2904938      total: 19.4s    remaining: 3s
866:      learn: 0.2904242      total: 19.4s    remaining: 2.98s
867:      learn: 0.2903578      total: 19.4s    remaining: 2.96s
```

In [21]:

```
1  # #####
2  # logging.warning("Begin Model.")
3  # result = {} # 保存一些结果
4  # # lightGBM 模型
5  # model_name = 'xgboost'
6  # version = 'v0.4_test'
7  # import xgboost as xgb
8
9  # # 5.1 模型定义
10 # model = xgb.XGBClassifier(
11 #     verbose=1
12 # )
13
14 # # 5.2 模型训练
15 # import time
16 # time_start=time.time()
17 # model.fit(X_train, Y_train)
18 # time_end=time.time()
19 # print('time cost',time_end-time_start,'s')
20 # # 5.3 模型评估
21 # result['score_train'] = model.score(X_train, Y_train)
22 # result['score_test'] = model.score(X_test, Y_test)
23 # # 5.4 模型保存
24 # joblib.dump(model, f'./model/%s_%s_[%f, %f, %f].model' %
25 #             (model_name, version,
26 #              result['score_test'], result['score_train'],
27 #              test_size))
28 # # 5.5 结果保存
29 # print(result)
30 # #####
```

WARNING:root:Begin Model.

/home/moegwjawel/atec_project/train/lmf/lib/python3.6/site-packages/xgboost/sklearn.py:1224: UserWarning: The use of label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove this warning, do the following: 1) Pass option use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
warnings.warn(label_encoder_deprecation_msg, UserWarning)

[19:10:54] WARNING: ../src/learner.cc:576:
Parameters: { "verbose" } might not be used.

This could be a false alarm, with some parameters getting used by language bindings but then being mistakenly passed down to XGBoost core, or some parameter actually being used

but getting flagged wrongly here. Please open an issue if you find any such cases.

[19:10:54] WARNING: ../src/learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
time cost 12.514598608016968 s
{'score_train': 0.9287125261445504, 'score_test': 0.8350180085976531}

In []:

```
1 #####  
2 # 下面是 rank 代码  
3 #####
```


In [22]:

```
1 # 2. 数据预处理
2 pd_data_rank = pd_raw_data.copy(deep=True)
3 print(pd_data_rank.shape)
4 # 2.1 去除 label 列
5 pd_data_rank = pd_data_rank.drop(columns=['label'], axis=1)
6 print(pd_data_rank.shape)
7 # 2.2 去除文字列
8 pd_data_rank = pd_data_rank.drop(columns=['memo_polish'], axis=1)
9 print(pd_data_rank.shape)
10 # 2.3 去除离群点
11 min_max_map = {}
12 file = open('min_max_map.txt', 'r')
13 for line in file.readlines():
14     line = line.strip()
15     line = line.split(' ')
16     min_max_map[line[0]] = (float(line[1]), float(line[2]))
17 file.close()
18 columns = pd_data_rank.columns.values.tolist()[1:] # 去除 id 的所有列
19 for col in columns:
20     (d_min, d_max) = min_max_map[col]
21     print(col, d_min, d_max)
22     pd_data_rank[pd_data_rank[col] > d_max] = np.nan
23     pd_data_rank[pd_data_rank[col] < d_min] = np.nan
24 print(pd_data_rank.shape)
25 # 2.4 去除取值只有一个值的列
26 drop_list = ['x2', 'x8', 'x13', 'x18', 'x25', 'x32', 'x33', 'x42', 'x48', 'x55',
27              'x79', 'x91', 'x96', 'x97', 'x100', 'x107', 'x110', 'x111', 'x113',
28              'x145', 'x152', 'x155', 'x178', 'x184', 'x198', 'x202', 'x207', 'x2',
29              'x232', 'x239', 'x248', 'x249', 'x258', 'x261', 'x264', 'x265', 'x2',
30              'x292', 'x300', 'x319', 'x325', 'x341', 'x343', 'x347', 'x351', 'x3',
31              'x389', 'x394', 'x400', 'x419', 'x427', 'x431', 'x436', 'x441', 'x4',
32              'x462', 'x475']
33 pd_data_rank = pd_data_rank.drop(columns=drop_list)
34 print(pd_data_rank.shape)
```

x402 -100.0 5.0
x403 -100.0 0.0
x404 -1.0 91656.79684362357
x405 -1.0 1.0
x406 0.0 2.0
x407 -1.0 0.05225940799998332
x408 -1.001 6.624499999998079
x409 0.0 23.0
x410 -1.001 5.0
x411 0.0 76251.99999994948
x412 -100.0 12.0
x413 -1.001 3.0
x414 -100.0 1.0
x415 0.0 1.0
x416 -1.001 17903.476899999645
x417 0.0 6.0
x418 0.0 129821.6106030297
x419 -1.002 -1.002
x420 -100.0 0.0
x421 -100.0 30.350000000000582

In [23]:

```
1 # 3. 数据转 numpy
2 np_data = pd_data_rank.values
3 # 去除索引/列
4 np_data = np_data[:, 1:]
5 # 分成 input 和 output, output 降维到一维
6 np_data_input = np_data[:, :np_data.shape[1]-1]
7 np_data_input = np_data_input.astype(np.float64)
8
9 print(np_data_input.shape)
10 print(np_data_input)
```

(43031, 398)

```
[[ 1.      -2.      0.      ... -100.      0.      1.      ]
 [ 1.      90.      0.      ... -100.      0.      1.      ]
 [ 1.      69.      0.      ... -100.      0.      7.      ]
 ...
 [ 1.      -2.      2.      ... -100.      0.      1.      ]
 [ 1.     179.      0.      ... -100.     60.42 -100.      ]
 [ 1.      -2.     -1.001 ... -100.      0.      7.      ]]
```

In [26]:

#15. 模型预测

```
output_predictions_path = './predictions.jsonl'
3
model_path_list = [
5 # ('./lightGBM_(0.010000, 0.858685, 0.853608).model', 1.0),
6 # ('./CatBoost_(0.010000, 0.876104, 0.845361).model', 1.0),
7 # ('./gbc_(0.010000, 0.830576, 0.820619).model', 1.0),
8 # ('./rfc_(0.010000, 0.999979, 0.830928).model', 1.0),
9 # ('./XGBoost_(0.010000, 0.902796, 0.837113).model', 1.0),
10 # ('./lightGBM_(0.100000, 0.858832, 0.848185).model', 1.0),
11 # ('./CatBoost_(0.100000, 0.877650, 0.845916).model', 1.0),
12 # ('./XGBoost_(0.100000, 0.905705, 0.840553).model', 1.0),
13 # ('./gbc_(0.100000, 0.831282, 0.836634).model', 1.0),
14 # ('./rfc_(0.100000, 0.999977, 0.840347).model', 1.0),
15
16 # ('./rfc.model', 0.003945), # 0.003945
17 # ('./lightGBM.model', 0.054045), # 0.054045
18 # ('./model/lightGBM_v0.1_test_[0.846638, 0.860215, 0.200000, LGBMClassifier(verb
19 ('./model/lightGBM_v0.4_test_[0.841641, 0.868783, 0.200000, LGBMClassifier(verb
20
pred_list = []
22
for model_path in model_path_list:
24 model = joblib.load(model_path[0])
25 model_score = model_path[1]
26 pred = model.predict_proba(np_data_input)
27 pred_list.append(
28     (pred, model_score)
29 )
30
#16. 保存结果
with open(output_predictions_path, 'w') as fp:
33 for i in range(np_data.shape[0]):
34     label = 0
35     score = 0
36     for pred in pred_list:
37         label += pred[0][i][1] * pred[1]
38         # print(pred[0][i][1], pred[1])
39         score += pred[1]
40     label /= score
41     # print("label:", label)
42     fp.write(f'{"id": "{i}", "label": {label}}\n')
43     logging.warning(f'{"id": "{i}", "label": {label}}\n')
44
```

WARNING:root:{"id": "2078", "label": 0.102900}

WARNING:root:{"id": "2079", "label": 0.773469}

WARNING:root:{"id": "2080", "label": 0.296777}

WARNING:root:{"id": "2081", "label": 0.013172}

WARNING:root:{"id": "2082", "label": 0.735827}

WARNING:root:{"id": "2083", "label": 0.154012}

WARNING:root:{"id": "2084", "label": 0.877750}

WARNING:root:{"id": "2085", "label": 0.451492}

```
WARNING:root:{ "id": "2085", "label": 0.703192},
```

```
WARNING:root:{ "id": "2086", "label": 0.703396}
```

```
WARNING:root:{ "id": "2087", "label": 0.627192}
```

In []:

1	
---	--