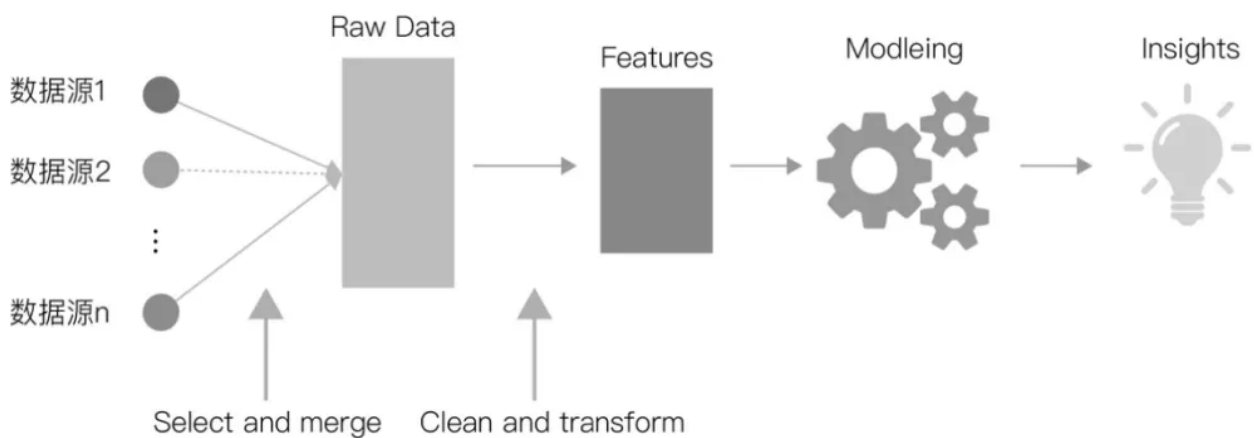


天池项目总结，特征工程了解一下！

阿里云天池 SAMshare 2020-12-22

来源：阿里云天池，案例：机器学习实践

业界广泛流传着这样一句话：“数据和特征决定了机器学习的上限，而模型和算法只是逼近这个上限而已”，由此可见特征工程在机器学习中的重要性，今天我们将通过《阿里云天池大赛赛题解析——机器学习篇》中的【天猫用户重复购买预测】案例来深入解析特征工程在实际商业场景中的应用。

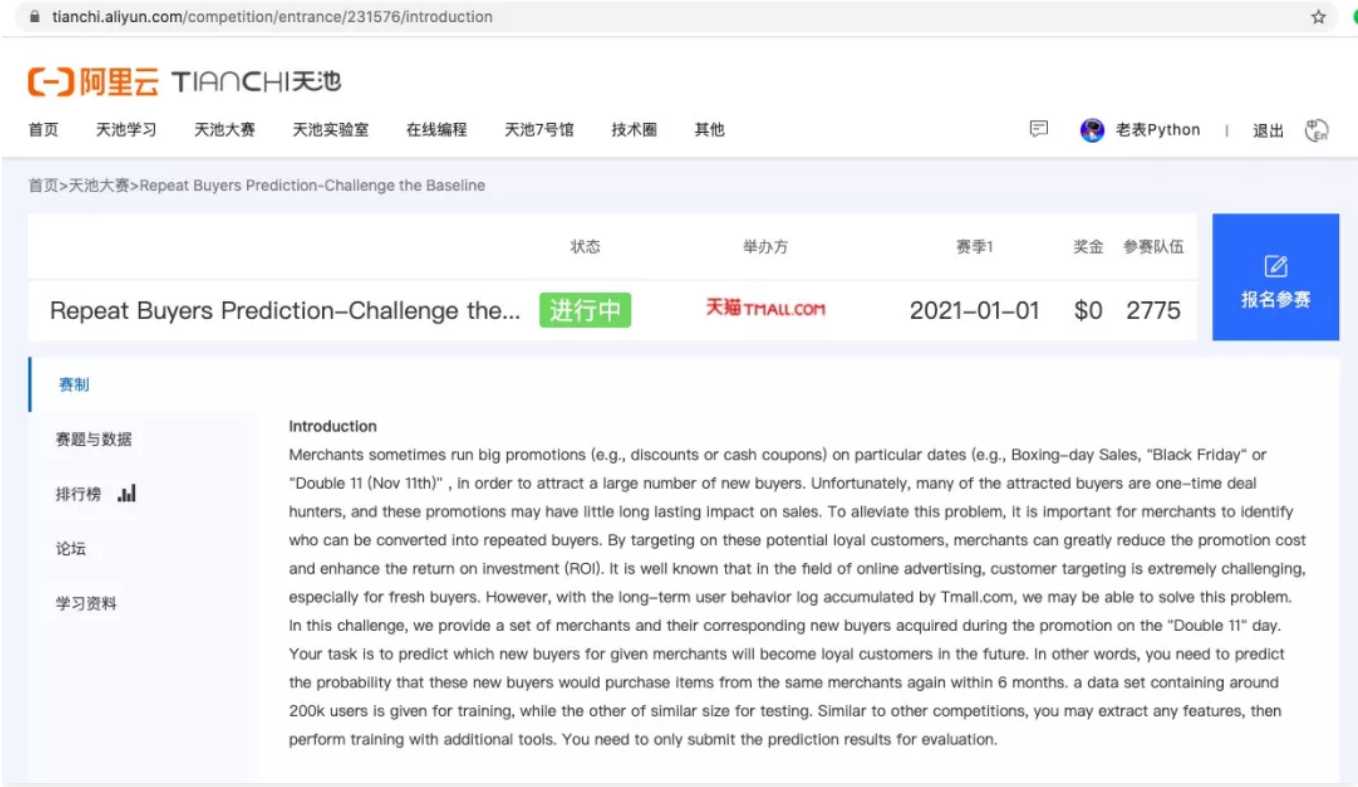


学习前须知

(1) 本文特征工程讲解部分参考自图书《阿里云天池大赛赛题解析——机器学习篇》中的第二个赛题：天猫用户重复购买预测。

(2) 本文相关数据可以在阿里云天池竞赛平台下载，数据地址：

<https://tianchi.aliyun.com/competition/entrance/231576/information>



一 数据集介绍

按照上面方法下载好数据集后，我们来看看具体数据含义。

data_format1	
• test_format1.csv	3.3MB
• train_format1.csv	3.5MB
• user_info_format1.csv	4.5MB
• user_log_format1.csv	1.91GB
data_format2	
• test_format2.csv	768.7MB
• train_format2.csv	767.4MB
sample_submission.csv	4.1MB

test_format1.csv和train_format1.csv里分别存放着测试数据和训练数据，测试数据集最后一个字段为prob，表示测试结果，训练数据集最后一个字段为label，训练数据各字段信息如下图所示：

Data Fields	Definition
user_id	顾客（用户）ID
merchant_id	商家（店铺）ID
label	取值范围：{0, 1}，1 表示重复购买，0 表示非重复购买，测试集数据的这个字段为 Null

训练数据集

user_log_format1.csv里存放着用户行为日志，字段信息如下图所示：

Data Fields	Definition
user_id	顾客（用户）ID
item_id	商品 ID
cat_id	商品类目 ID
seller_id	商家（店铺）ID merchantid
brand_id	品牌 ID
time_tamp	行为发生时间
action_type	取值范围：{0, 1, 2, 3}，0 表示点击，1 表示加入购物车，2 表示购买，3 表示收藏

用户行为日志数据

user_info_format1.csv里存放着用户行基本信息，字段信息如下图所示：

Data Fields	Definition
user_id	顾客（用户）ID
age_range	顾客（用户）年龄范围：1 表示<18；2 表示[18， 24]；3 表示[25， 29]；4 表示[30， 34]；5 表示[35， 39]；6 表示[40， 49]；7 and 8 表示>=50；0and Null 表示未知
gender	顾客性别：0 表示女性；1 表示男性；2and Null 表示未知

用户基本信息数据

二 特征构造

本赛题基于天猫电商数据，主要关心用户、店铺和商家这三个实体，所以特征构造上也以用户、店铺和商家为核心，可以分为以下几部分：



用户-店铺特征构造

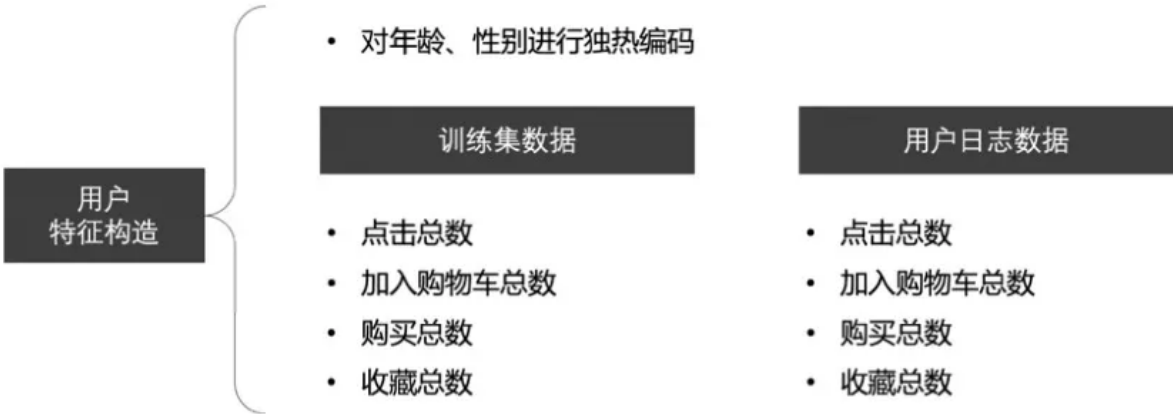


店铺特征构造

对店铺特征选取可以使用，如 Numpy 的 `corrcoef(x,y)`函数计算相关系数，保留相关系数小于0.9 的特征组合，具体内容如图 2-3。



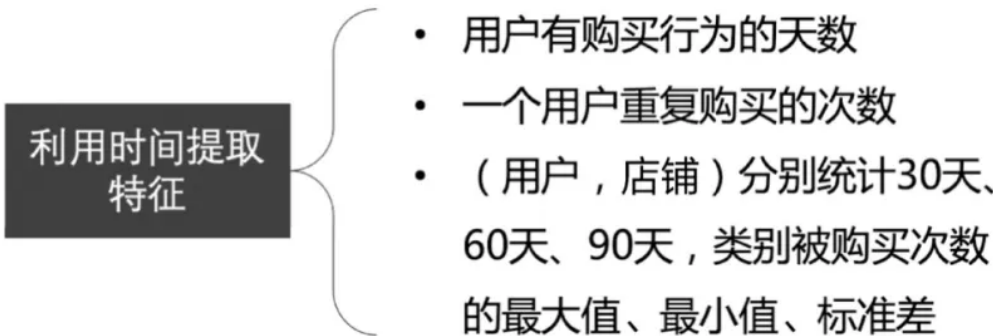
商家特征选取



用户特征构造

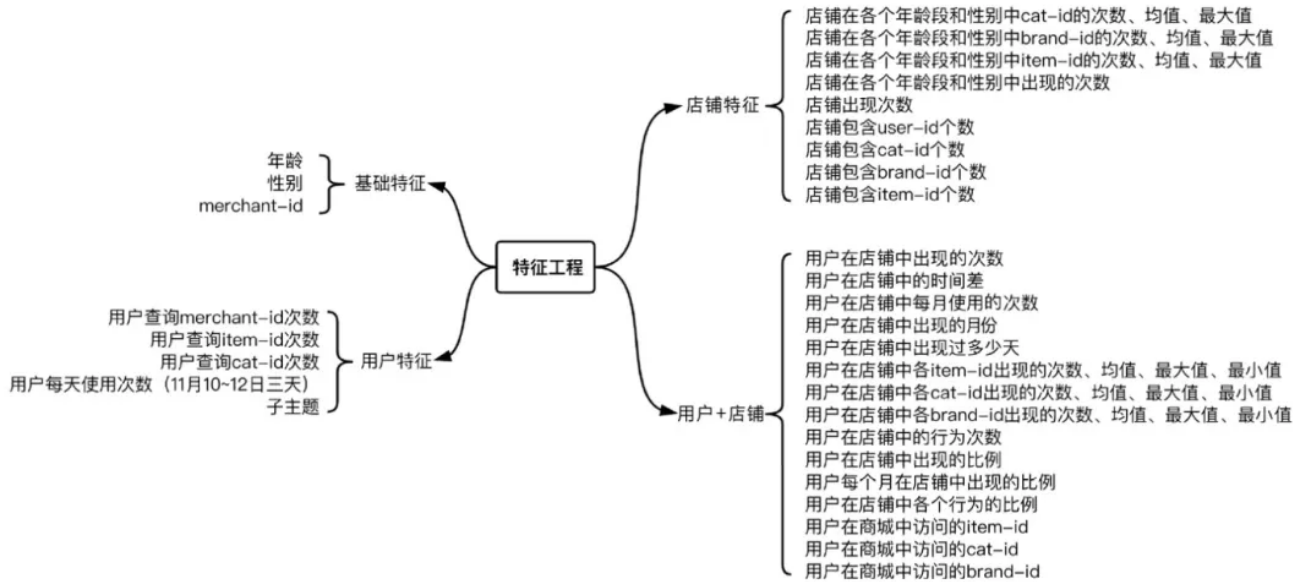


用户购买商品特征构造



利用时间提取特征

总结以上内容，特征主要基于基础特征、用户特征、店铺特征、用户+店铺四个方面，如下图所示：



特征总结

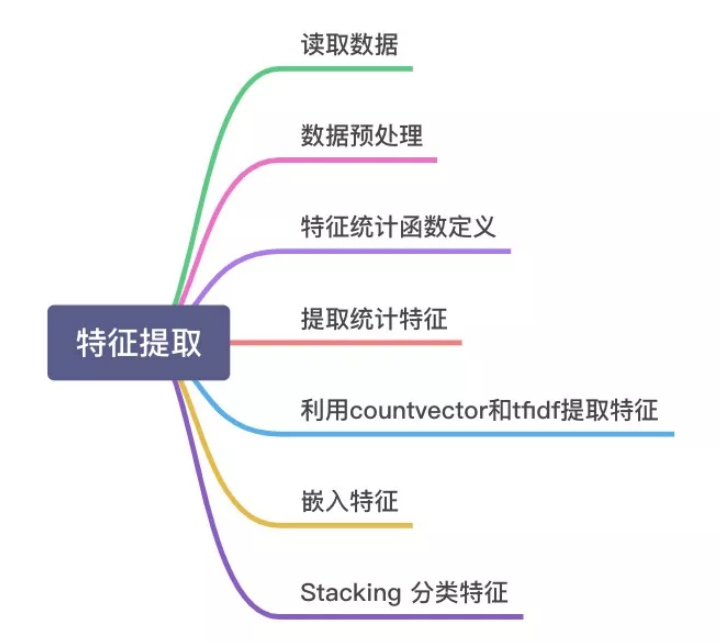
三 特征提取

首先我们导入需要的工具包，进行数据分析和特征提取。

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt import seaborn as sns
4 from scipy import stats
5 import gc
6 from collections import Counter import copy
7 import warnings warnings.filterwarnings("ignore")
8 %matplotlib inline
  
```

接下来我们将按以下步骤进行特征提取。



特征提取步骤

1 读取数据

直接调用Pandas的read_csv函数读取训练集和测试集及用户信息、用户日志数据。

```
1 test_data = pd.read_csv('./data_format1/test_format1.csv') train_data = pd.read_csv('./data_format1/train_format1.csv')
```

2 数据预处理

对数据内存进行压缩：

```

1 def reduce_mem_usage(df, verbose=True):
2     start_mem = df.memory_usage().sum() / 1024**2
3     numerics = ['int16', 'int32', 'int64', 'float16', 'float32', 'float64']
4     for col in df.columns: col_type = df[col].dtypes if col_type in numerics:
5         c_min = df[col].min()
6         c_max = df[col].max()
7         if str(col_type)[:3] == 'int':
8             if c_min > np.iinfo(np.int8).min and c_max < np.iinfo(np.int8).max:
9                 df[col] = df[col].astype(np.int8)
10            elif c_min > np.iinfo(np.int16).min and c_max < np.iinfo(np.int16).max:
11                df[col] = df[col].astype(np.int16)

```

```

12         df[col] = df[col].astype(np.int16)
13     elif c_min > np.iinfo(np.int32).min and c_max < np.iinfo(np.int32).max:
14         df[col] = df[col].astype(np.int32)
15     elif c_min > np.iinfo(np.int64).min and c_max < np.iinfo(np.int64).max:
16         df[col] = df[col].astype(np.int64)
17     else:
18         if c_min > np.finfo(np.float16).min and c_max < np.finfo(np.float16).max:
19             df[col] = df[col].astype(np.float16)
20         elif c_min > np.finfo(np.float32).min and c_max < np.finfo(np.float32).max:
21             df[col] = df[col].astype(np.float32)
22         else:
23             df[col] = df[col].astype(np.float64)
24     end_mem = df.memory_usage().sum() / 1024**2
25     print('Memory usage after optimization is: {:.2f} MB'.format(end_mem))

```

首先测试数据添加到训练数据后, 然后将用户基本信息合并到训练数据左边, 并删除不需要的变量, 释放内存。

```

1 all_data = train_data.append(test_data)
2 all_data = all_data.merge(user_info, on='user_id', how='left')
3 del train_data, test_data, user_info
4 gc.collect()

```

将用户日志数据各字段合并成一个新的字段item_id, 并将其插入到用户信息数据之后。

```

1 # 用户日志数据按时间排序
2 user_log = user_log.sort_values(['user_id', 'time_stamp'])
3 # 合并用户日志数据各字段, 新字段名为item_id
4 list_join_func = lambda x: " ".join([str(i) for i in x])
5 agg_dict = {
6     'item_id': list_join_func,
7     'cat_id': list_join_func,

```



```

8     'seller_id': list_join_func,
9     'brand_id': list_join_func,
10    'time_stamp': list_join_func,
11    'action_type': list_join_func
12 }
13 rename_dict = {
14     'item_id': 'item_path',
15     'cat_id': 'cat_path',
16     'seller_id': 'seller_path',
17     'brand_id': 'brand_path',
18     'time_stamp': 'time_stamp_path',
19     'action_type': 'action_type_path'
20 }
21
22 def merge_list(df_ID, join_columns, df_data, agg_dict, rename_dict):
23     df_data = df_data.groupby(join_columns).agg(agg_dict).reset_index().rename(
24         columns=rename_dict)
25     df_ID = df_ID.merge(df_data, on=join_columns, how="left")
26     return df_ID
27 all_data = merge_list(all_data, 'user_id', user_log, agg_dict, rename_dict)
28 del user_log
29 gc.collect()

```

3 特征统计函数定义

基于之前的特征构造图, 我们提前编写一些统计相关函数, 依次有: 数据总数、数据唯一值总数、数据最大值、数据最小值、数据标准差、数据中top N数据以及数据中top N数据的总数。

```

1 def cnt_(x):
2     try:
3         return len(x.split(' '))
4     except:
5         return -1
6
7 def nunique_(x):
8     try:
9         return len(set(x.split(' ')))

```

```
10     except:
11         return -1
12
13 def max_(x):
14     try:
15         return np.max([float(i) for i in x.split(' ')])
16     except:
17         return -1
18
19 def min_(x):
20     try:
21         return np.min([float(i) for i in x.split(' ')])
22     except:
23         return -1
24
25 def std_(x):
26     try:
27         return np.std([float(i) for i in x.split(' ')])
28     except:
29         return -1
30
31 def most_n(x, n):
32     try:
33         return Counter(x.split(' ')).most_common(n)[n-1][0]
34     except:
35         return -1
36
37 def most_n_cnt(x, n):
38     try:
39         return Counter(x.split(' ')).most_common(n)[n-1][1]
40     except:
41         return -1
```

基于上面编写的基本统计方法, 我们可以针对数据进行特征统计。

```
1 def user_cnt(df_data, single_col, name):
2     df_data[name] = df_data[single_col].apply(cnt_)
3
```

```
4     return df_data
5
6 def user_nunique(df_data, single_col, name):
7     df_data[name] = df_data[single_col].apply(nunique_)
8     return df_data
9
10 def user_max(df_data, single_col, name):
11     df_data[name] = df_data[single_col].apply(max_)
12     return df_data
13
14 def user_min(df_data, single_col, name):
15     df_data[name] = df_data[single_col].apply(min_)
16     return df_data
17
18 def user_std(df_data, single_col, name):
19     df_data[name] = df_data[single_col].apply(std_)
20     return df_data
21
22 def user_most_n(df_data, single_col, name, n=1):
23     func = lambda x: most_n(x, n)
24     df_data[name] = df_data[single_col].apply(func)
25     return df_data
26
27 def user_most_n_cnt(df_data, single_col, name, n=1):
28     func = lambda x: most_n_cnt(x, n)
29     df_data[name] = df_data[single_col].apply(func)
30     return df_data
```

4 提取统计特征

基于上一步中编写的用户数据统计函数, 以店铺特征统计为例, 统计与店铺特点有关的特征, 如店铺、商品、品牌等。

```
1 # 取2000条数据举例
2 all_data_test = all_data.head(2000)
3 # 总次数
4 all_data_test = user_cnt(all_data_test, 'seller_path', 'user_cnt')
```

```

5 # 不同店铺个数
6 all_data_test = user_nunique(all_data_test, 'seller_path', 'seller_nunique ')
7 # 不同品类个数
8 all_data_test = user_nunique(all_data_test, 'cat_path', 'cat_nunique')
9 # 不同品牌个数
10 all_data_test = user_nunique(all_data_test, 'brand_path',
11                             'brand_nunique')
12 # 不同商品个数
13 all_data_test = user_nunique(all_data_test, 'item_path', 'item_nunique')
14 # 活跃天数
15 all_data_test = user_nunique(all_data_test, 'time_stamp_path',
16                             'time_stamp _nunique')
17 # 不同用户行为种数
18 all_data_test = user_nunique(all_data_test, 'action_type_path',
19                             'action_ty pe_nunique')

```

此外还可以统计用户最喜欢的店铺、最喜欢的类目、最喜欢的品牌、最长见的行为动作等数据。

```

1 # 用户最喜欢的店铺
2 all_data_test = user_most_n(all_data_test, 'seller_path', 'seller_most_1', n=1)
3 # 最喜欢的类目
4 all_data_test = user_most_n(all_data_test, 'cat_path', 'cat_most_1', n=1)
5 # 最喜欢的品牌
6 all_data_test = user_most_n(all_data_test, 'brand_path', 'brand_most_1', n= 1)
7 # 最常见的行为动作
8 all_data_test = user_most_n(all_data_test, 'action_type_path', 'action_type _1

```

5 利用countvector和tfidf提取特征

CountVectorizer与TfidfVectorizer是Scikit-learn的两个特征数值计算的类, 接下来我们将结合两者进行特征提取。

```

1 from sklearn.feature_extraction.text import CountVectorizer
2

```

```
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.feature_extraction.text import ENGLISH_STOP_WORDS
5 from scipy import sparse
6 tfidfVec = TfidfVectorizer(stop_words=ENGLISH_STOP_WORDS,
7                             ngram_range=(1, 1),
8                             max_features=100)
9 columns_list = ['seller_path']
10 for i, col in enumerate(columns_list):
11     tfidfVec.fit(all_data_test[col])
12     data_ = tfidfVec.transform(all_data_test[col])
13     if i == 0:
14         data_cat = data_
15     else:
16         data_cat = sparse.hstack((data_cat, data_))
```

6 嵌入特征

词嵌入是一类将词向量化的模型的统称, 核心思想是将每个词都映射到低维空间(K 为50~300)上的一个稠密向量。

```
1 import gensim
2
3 model = gensim.models.Word2Vec(
4     all_data_test['seller_path'].apply(lambda x: x.split(' ')),
5     size=100,
6     window=5,
7     min_count=5,
8     workers=4)
9
10 def mean_w2v_(x, model, size=100):
11     try:
12         i = 0
13         for word in x.split(' '):
14             if word in model.wv.vocab:
15                 i += 1
16             if i == 1:
17                 vec = np.zeros(size)
```

```

18         vec += model.wv[word]
19     return vec / i
20 except:
21     return np.zeros(size)
22
23 def get_mean_w2v(df_data, columns, model, size):
24     data_array = []
25     for index, row in df_data.iterrows():
26         w2v = mean_w2v_(row[columns], model, size)
27         data_array.append(w2v)
28     return pd.DataFrame(data_array)
29
30 df_embeeding = get_mean_w2v(all_data_test, 'seller_path', model, 100)
31 df_embeeding.columns = ['embeeding_' + str(i) for i in df_embeeding.columns]

```

7 Stacking 分类特征

以使用 lgb 和 xgb 分类模型构造 Stacking 特征为例子, 实现方式如下:

```

1  # 1、使用 5 折交叉验证
2  from sklearn.model_selection import StratifiedKFold, KFold
3  folds = 5
4  seed = 1
5  kf = KFold(n_splits=5, shuffle=True, random_state=0)
6  # 2、选择 lgb 和 xgb 分类模型作为基模型
7  clf_list = [lgb_clf, xgb_clf]
8  clf_list_col = ['lgb_clf', 'xgb_clf']
9  # 3、获取 Stacking 特征
10 clf_list = clf_list
11 column_list = []
12 train_data_list=[]
13 test_data_list=[]
14 for clf in clf_list:
15     train_data,test_data,clf_name=clf(x_train, y_train, x_valid, kf, label_sp
16     train_data_list.append(train_data)
17     test_data_list.append(test_data)
18 train_stacking = np.concatenate(train_data_list, axis=1)

```

```
19 test_stacking = np.concatenate(test_data_list, axis=1)
```

运行结果:

```
1 [1] valid_0's multi_logloss: 0.240875
2 Training until validation scores don't improve for 100 rounds.
3 [2] valid_0's multi_logloss: 0.240675
4 [226] train-mlogloss:0.123211 eval-mlogloss:0.226966
5 Stopping. Best iteration:
6 [126] train-mlogloss:0.172219 eval-mlogloss:0.218029
7 xgb now score is: [2.4208301225770263, 2.2433633135072886, 2.51909203146584 34
8 xgb_score_list: [2.4208301225770263, 2.2433633135072886, 2.5190920314658434, 2
9 xgb_score_mean: 2.4506746084485203
```

SAMshare

置顶我哟, 不要走宝啦



"干货学习, 分享在看三连↓"

喜欢此内容的人还喜欢

三万字, Spark学习笔记

数据社

2021年人社部更新《国家职业资格目录(专业技术人员职业资格)》

川杨学堂