

## kaggle系列（3）：Rental Listing Inquiries（二）：XGBoost

📅 2017-06-13 | 📁 Kaggle | 📄 91

上一节我们对数据集进行了初步的探索，并将其可视化，对数据有了初步的了解。这样我们有了之前数据探索的基础之后，就有了对其建模的基础feature，结合目标变量，即可进行模型训练了。我们使用交叉验证的方法来判断线下的实验结果，也就是把训练集分成两部分，一部分是训练集，用来训练分类器，另一部分是验证集，用来计算损失评估模型的好坏。

在Kaggle的希格斯子信号识别竞赛中，XGBoost因为出众的效率与较高的预测准确度在比赛论坛中引起了参赛选手的广泛关注，在1700多支队伍的激烈竞争中占有一席之地。随着它在Kaggle社区知名度的提高，最近也有队伍借助XGBoost在比赛中夺得第一。其次，因为它的效果好，计算复杂度不高，也在工业界中有大量的应用。

今天，我们就先来跑一个XGBoost版的Base Model。先回顾一下XGBoost的原理吧：机器学习算法系列（8）：XgBoost

### 一、准备工作

首先我们导入需要的包：

```
1  import os
2  import sys
3  import operator
4  import numpy as np
5  import pandas as pd
6  from scipy import sparse
7  import xgboost as xgb
8  from sklearn import model_selection, preprocessing, ensemble
9  from sklearn.metrics import log_loss
10 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
```

其中一些包的用途会在之后具体用到的时候进行讲解。

导入我们的数据：

```
1 data_path = '../data/'
2 train_file = data_path + "train.json"
3 test_file = data_path + "test.json"
4 train_df = pd.read_json(train_file)
5 test_df = pd.read_json(test_file)
6 print train_df.shape
7 print test_df.shape
8
9 (49352, 15)
10 (74659, 14)
```

查看一下前两行：

```
1 train_df.head(2)
```

bathrooms	bedrooms	building_id	created	description	display_address	features	interest_level	
10	1.5	3	53a5b119ba8f7b61d4e010512e0dfc85	2016-06-24 07:54:24	A Brand New 3 Bedroom 1.5 bath ApartmentEnjoy ...	Metropolitan Avenue	[]	medium
10000	1.0	2	c5c8a357cba207596b04d1afd1e4f130	2016-06-12 12:19:27		Columbus Avenue	[Doorman, Elevator, Fitness Center, Cats Allow...	low

latitude	listing_id	longitude	manager_id	photos	price	street_address
40.7145	7211212	-73.9425	5ba989232d0489da1b5f2c45f6688adc	[https://photos.renthop.com/2/7211212_1ed4542e...	3000	792 Metropolitan Avenue
40.7947	7150865	-73.9667	7533621a882f71e25173b27e3139d83d	[https://photos.renthop.com/2/7150865_be3306c5...	5465	808 Columbus Avenue

## 二、特征构建

我们不需要对数值型数据进行任何的预处理，所以首先建立一个数值型特征的列表，纳入 features\_to\_use

```
1 features_to_use = ["bathrooms", "bedrooms", "latitude", "longitude", "price"]
```

现在让我们根据已有的一些特征来构建一些新的特征：

```
1 # 照片数量(num_photos)
2 train_df['num_photos']=train_df['photos'].apply(len)
3 test_df['num_photos']=train_df['photos'].apply(len)
4
5 # 特征数量
6 train_df['num_features']=train_df['features'].apply(len)
7 test_df['num_features']=test_df['features'].apply(len)
8
9 # 描述词汇数量
10 train_df['num_description_words'] = train_df['description'].apply(lambda x: len(x.split))
11 test_df['num_description_words'] = test_df['description'].apply(lambda x: len(x.split("
12
13 #把创建的时间分解为多个特征
14 train_df['created']=pd.to_datetime(train_df['created'])
15 test_df['created']=pd.to_datetime(test_df['created'])
16
17 #让我们从时间中分解出一些特征，比如年，月，日，时
18 #年
19 train_df['created_year'] = train_df['created'].dt.year
20 test_df['created_year'] = test_df['created'].dt.year
21 #月
22 train_df['created_month'] = train_df['created'].dt.month
23 test_df['created_month'] = test_df['created'].dt.month
24 #日
25 train_df['created_day'] = train_df['created'].dt.day
26 test_df['created_day'] = test_df['created'].dt.day
27 #时
28 train_df['created_hour'] = train_df['created'].dt.hour
29 test_df['created_hour'] = test_df['created'].dt.hour
30
31 #把这些特征都放到所需特征列表中（上面已经创建，并加入了数值型特征）
32 features_to_use.extend(["num_photos", "num_features", "num_description_words", "created_ye
```

我们有四个分类型的特征：

- display\_address
- manager\_id
- building\_id
- street\_address

可以对它们分别进行特征编码：

```
1 categorical = ["display_address", "manager_id", 'building_id', "street_address"]
2 for f in categorical:
3     if train_df[f].dtype == 'object':
4         lbl = preprocessing.LabelEncoder()
5         lbl.fit(list(train_df[f].values)+list(test_df[f].values))
6         train_df[f] = lbl.transform(list(train_df[f].values))
7         test_df[f] = lbl.transform(list(test_df[f].values))
8         features_to_use.append(f)
```

还有一些字符串类型的特征，可以先把它们合并起来

```
1 train_df["features"] = train_df["features"].apply(lambda x: " ".join(["_".join(i.split("
2 print train_df['features'].head(2)
3 test_df['features'] = test_df["features"].apply(lambda x: " ".join(["_".join(i.split(" "
4 print test_df['features'].head(2)
```

得到的字符串结果如下：

10000 Doorman Elevator Fitness\_Center Cats\_Allowed D...

100004 Laundry\_In\_Building Dishwasher Hardwood\_Floors...

然后CountVectorizer类来计算TF-IDF权重

```
1 tfidf = CountVectorizer(stop_words ="english",max_features=200)
2 tr_sparse = tfidf.fit_transform(train_df["features"])
3 te_sparse = tfidf.transform(test_df["features"])
```

这里我们需要提一点，对数据集进行特征变换时，必须同时对训练集和测试集进行操作。现在把这些处理过的特征放到一个集合中（横向合并）

```
1 train_X = sparse.hstack([train_df[features_to_use],tr_sparse]).tocsr()
2 test_X = sparse.hstack([test_df[features_to_use],te_sparse]).tocsr()
```

然后把目标变量转换为0、1、2，如下

```

1 target_num_map = {'high':0 , 'medium':1 , 'low':2}
2 train_y = np.array(train_df['interest_level'].apply(lambda x: target_num_map[x]))
3 print train_X.shape,test_X.shape
4
5 (49352, 217) (74659, 217)

```

可以看到，经过上面一系列的变量构造之后，其数量已经达到了217个。

接下来就可以进行建模啦。

### 三、XGB建模

先写一个通用的XGB模型的函数：

```

1 def runXGB(train_X,train_y,test_X,test_y=None,feature_names=None,seed_val=0,num_rounds=
2     #参数设定
3     param = {}
4     param['objective'] = 'multi:softprob'#多分类、输出概率值
5     param['eta'] = 0.1#学习率
6     param['max_depth'] = 6#最大深度，越大越容易过拟合
7     param['silent'] = 1#打印提示信息
8     param['num_class'] = 3#三个类别
9     param['eval_metric'] = 'mlogloss'#对数损失
10    param['min_child_weight']=1#停止条件，这个参数非常影响结果，控制叶子节点中二阶导的和的最
11    param['subsample'] = 0.7#随机采样训练样本
12    param['colsample_bytree'] = 0.7# 生成树时进行的列采样
13    param['seed'] = seed_val#随机数种子
14    num_rounds = num_rounds#迭代次数
15
16    plst = list(param.items())
17    xgtrain = xgb.DMatrix(train_X,label=train_y)
18
19    if test_y is not None:
20        xgtest = xgb.DMatrix(test_X,label=test_y)
21        watchlist = [(xgtrain,'train'),(xgtest,'test')]
22        model = xgb.train(plst,xgtrain,num_rounds,watchlist,early_stopping_rounds=20)
23        # early_stopping_rounds 当设置的迭代次数较大时，early_stopping_rounds 可在一定的迭代
24    else:
25        xgtest = xgb.DMatrix(test_X)
26        model = xgb.train(plst,xgtrain,num_rounds)
27    pred_test_y = model.predict(xgtest)
28    return pred_test_y,model

```

函数返回的是预测值和模型。

5折交叉验证将训练集划分为五份，其中的一份作为验证集。

```
1 cv_scores = []
2 kf = model_selection.KFold(n_splits=5,shuffle=True,random_state=2016)
3
4 for dev_index,val_index in kf.split(range(train_X.shape[0])):
5     dev_X,val_X = train_X[dev_index:],train_X[val_index,:]
6     dev_y,val_y = train_y[dev_index],train_y[val_index]
7     pred,model = runXGB(dev_X,dev_y,val_X,val_y)
8     cv_scores.append(log_loss(val_y,preds))
9     print cv_scores
10    break
```

结果如下：

```
1  [0]      train-mlogloss:1.04135  test-mlogloss:1.04229
2  Multiple eval metrics have been passed: 'test-mlogloss' will be used for early stopping
3
4  Will train until test-mlogloss hasn't improved in 20 rounds.
5  [1]      train-mlogloss:0.989004 test-mlogloss:0.99087
6  [2]      train-mlogloss:0.944233 test-mlogloss:0.947047
7  [3]      train-mlogloss:0.90536  test-mlogloss:0.908933
8  [4]      train-mlogloss:0.872054 test-mlogloss:0.876526
9  [5]      train-mlogloss:0.841783 test-mlogloss:0.847383
10 [6]      train-mlogloss:0.815921 test-mlogloss:0.822307
11 [7]      train-mlogloss:0.793337 test-mlogloss:0.800476
12 [8]      train-mlogloss:0.773562 test-mlogloss:0.781413
13 [9]      train-mlogloss:0.754927 test-mlogloss:0.76381
14 [10]     train-mlogloss:0.738299 test-mlogloss:0.747959
15
16 .....
17 .....
18 [367]    train-mlogloss:0.348196 test-mlogloss:0.548011
19 [368]    train-mlogloss:0.347768 test-mlogloss:0.547992
20 [369]    train-mlogloss:0.347303 test-mlogloss:0.548021
21 [370]    train-mlogloss:0.346807 test-mlogloss:0.548065
22 [371]    train-mlogloss:0.346514 test-mlogloss:0.548079
23 [372]    train-mlogloss:0.34615  test-mlogloss:0.548097
24 [373]    train-mlogloss:0.345859 test-mlogloss:0.548111
25 [374]    train-mlogloss:0.345377 test-mlogloss:0.548081
26 [375]    train-mlogloss:0.344961 test-mlogloss:0.548068
27 [376]    train-mlogloss:0.344493 test-mlogloss:0.548024
28 [377]    train-mlogloss:0.344086 test-mlogloss:0.547975
29 Stopping. Best iteration:
30 [357]    train-mlogloss:0.352182 test-mlogloss:0.547867
```

迭代357次之后，在训练集上的对数损失为0.352182，在验证集上的损失为0.5478。

然后在对测试集进行预测：

```
1 preds,model=runXGB(train_X,train_y,test_X,num_rounds=400)
```

把结果按照比赛规定的格式写入csv文件：

```
1 out_df = pd.DataFrame(preds)
2 out_df.columns = ["high", "medium", "low"]
3 out_df["listing_id"] = test_df.listing_id.values
4 out_df.to_csv("xgb_starter2.csv", index=False)
```

看一下最后的结果：

	high	medium	low	listing_id
0	0.015237	0.198407	0.786357	7142618
1	0.014637	0.049824	0.935538	7210040
2	0.002438	0.064852	0.932710	7103890
3	0.021119	0.480018	0.498863	7143442
4	0.008035	0.068569	0.923395	6860601
5	0.000116	0.019608	0.980276	6840081

提交到kaggle上，这样我们整个建模的过程就完成了。

接下来两节中，我们重点讲一讲关于XGBoost的调参经验以及使用SK-learn计算TF-IDF。

感谢打赏，您的支持将鼓励我继续创作！



# kaggle    # XGBoost    # 特征工程