

代码解读

```
1 # import packages
2 import time, math, os
3 from tqdm import tqdm
4 import gc
5 import pickle
6 import random
7 from datetime import datetime
8 from operator import itemgetter
9 import numpy as np
10 import pandas as pd
11 import warnings
12 from collections import defaultdict
13 import collections
14 warnings.filterwarnings('ignore')
```

operator模块中的itemgetter()函数!

```
1 a = [1,2,3]
2 >>> b=operator.itemgetter(1)      //定义函数b，获取对象的第1个域的值
3 >>> b(a)
4 2
5 >>> b=operator.itemgetter(1,0)    //定义函数b，获取对象的第1个域和第0个的值
6 >>> b(a)
7 (2, 1)
```

```
1 # 全量训练集
2 all_click_df = get_all_click_df(data_path, offline=False)
```

Pandas之drop_duplicates: 去除重复项

DataFrame.drop_duplicates(subset=None, keep='first', inplace=False)

- subset : column label or sequence of labels, optional
用来指定特定的列，默认所有列
- keep : {'first', 'last', False}, default 'first'
删除重复项并保留第一次出现的项
- inplace : boolean, default False
是直接原来数据上修改还是保留一个副本

```
In [22]: import pandas as pd
data=pd.DataFrame({'A':[1,1,2,2], 'B':['a','b','a','b']})
data.drop_duplicates('B', 'first', inplace = True)
data
```

```
Out[22]:
```

	A	B
0	1	a
1	1	b

<http://blog.csdn.net/u010665216>

此地无银三百两

```
1 def get_all_click_df(data_path='./data_raw/', offline=True):
2     if offline:
3         all_click = pd.read_csv(data_path + 'train_click_log.csv')
4     else:
5         trn_click = pd.read_csv(data_path + 'train_click_log.csv')
6         tst_click = pd.read_csv(data_path + 'testA_click_log.csv')
7
8         all_click = trn_click.append(tst_click)
9
10    all_click = all_click.drop_duplicates(['user_id', 'click_article_id',
11    'click_timestamp'])
12    # 笔者在后面才发现前面有问题 送命题请注意
13    # ['user_id', 'click_article_id', 'click_timestamp']要求三者参数都一样才能达到去重效果
14    return all_click
```

```
1 i2i_sim = itemcf_sim(all_click_df)
```

```
1 def itemcf_sim(df):
2     """
3     文章与文章之间的相似性矩阵计算
4     :param df: 数据表
5     :item_created_time_dict: 文章创建时间的字典
6     return : 文章与文章的相似性矩阵
7     思路: 基于物品的协同过滤(详细请参考上一期推荐系统基础的组队学习), 在多路召回部分会加上关联规则的召回策略
```

```

8     """
9     # df是合并起来的数据集
10    user_item_time_dict = get_user_item_time(df)
11    # get_user_item_time(df) 见下一块代码
12
13    # 计算物品相似度
14    i2i_sim = {}
15    item_cnt = defaultdict(int) # 初始化为0
16    # tqdm显示for循环完成的进度
17    for user, item_time_list in tqdm(user_item_time_dict.items()): # {user1:
        [(item1, time1), (item2, time2)..]...}
18        # 在基于商品的协同过滤优化的时候可以考虑时间因素
19        for i, i_click_time in item_time_list: # i:item    i_click_time:time
20            item_cnt[i] += 1 # 统计物品i被不同的人买了几次
21            i2i_sim.setdefault(i, {}) # i2i_sim空字典 {user1:{},user2:{},...}
22            for j, j_click_time in item_time_list: # j=item    j_click_time=time
23                if(i == j):
24                    continue
25                i2i_sim[i].setdefault(j, 0) # { item1:{item2:0,item3:0,...}, item2:
        {item1:0,item3:0,...},...}
26
27                i2i_sim[i][j] += 1 / math.log(len(item_time_list) + 1) # 对活跃用户的
        惩罚 另一种计算相似度公式
28
29    i2i_sim_ = i2i_sim.copy() # copy一份字典
30    for i, related_items in i2i_sim.items(): # i=item1    related_items={item2:
        i2i_sim[i][j],item3: i2i_sim[i][j],...}
31        for j, wij in related_items.items(): # j=item2    wij=i2i_sim[i][j]
32            i2i_sim_[i][j] = wij / math.sqrt(item_cnt[i] * item_cnt[j]) #
        item_cnt[i]物品i被统计次数 item_cnt[j]物品j被统计次数 i!=j
33
34    # 将得到的相似性矩阵保存到本地
35    pickle.dump(i2i_sim_, open(save_path + 'itemcf_i2i_sim.pkl', 'wb')) # 它可以将对
        象转换为一种可以传输或存储的格式。
36
37    return i2i_sim_

```

tqdm

获取 用户 - 文章 - 点击时间字典

```

1  # 根据点击时间获取用户的点击文章序列    {user1: [(item1, time1), (item2, time2)..]...}
2  def get_user_item_time(click_df):
3      # click_df参数是pandas中pd合并后的数据集
4      click_df = click_df.sort_values('click_timestamp')
5      # 就是按照表头click_timestamp大小排序
6      # 送命题 一定要有click_df =
7      # 否则你print(click_df)和print(click_df.sort_values('click_timestamp'))输出不一样

```

```

8
9     def make_item_time_pair(df):
10         # 拉锁函数
11         return list(zip(df['click_article_id'], df['click_timestamp']))
12
13     user_item_time_df = click_df.groupby('user_id')['click_article_id',
14     'click_timestamp'].apply(
15         lambda x: make_item_time_pair(x)) /
16         .reset_index().rename(columns={0: 'item_time_list'})
17     user_item_time_dict = dict(zip(user_item_time_df['user_id'],
18     user_item_time_df['item_time_list']))
19
20     return user_item_time_dict

```

拉锁函数超链接 [Pandas: apply方法与lambda、groupby结合、apply多参数传递](#)

Pandas之sort_values isin使用技巧

```

1 import pandas as pd
2 df = pd.DataFrame([[1,2,3],[2,3,4],[2,4,3],[1,3,7]],
3     index = ['one','two','three','four'],columns = ['A','B','C'])
4 print df
5 #      A  B  C
6 # one  1  2  3
7 # two  2  3  4
8 # three 2  4  3
9 # four  1  3  7
10 df.sort_values(by=['A','B'],ascending=[0,1],inplace=True)
11 print df
12 #      A  B  C
13 # two  2  3  4
14 # three 2  4  3
15 # one  1  2  3
16 # four  1  3  7
17

```

```

1 # 定义
2 user_recall_items_dict = collections.defaultdict(dict)
3 # <class 'dict'> dict
4 # defaultdict(<class 'dict'>, {}) user_recall_items_dict
5
6 # 获取 用户 - 文章 - 点击时间的字典
7 user_item_time_dict = get_user_item_time(all_click_df) # 参数是合并后的数据集
8

```

```

1 # 根据点击时间获取用户的点击文章序列 {user1: [(item1, time1), (item2, time2)..]...}
2 def get_user_item_time(click_df):
3
4     click_df = click_df.sort_values('click_timestamp') # 排序
5
6     def make_item_time_pair(df):
7         return list(zip(df['click_article_id'], df['click_timestamp']))
8
9     user_item_time_df = click_df.groupby('user_id')['click_article_id',
10     'click_timestamp'].apply(lambda x: make_item_time_pair(x)) /
11
12     .reset_index().rename(columns={0: 'item_time_list'})
13     user_item_time_dict = dict(zip(user_item_time_df['user_id'],
14     user_item_time_df['item_time_list']))
15
16     return user_item_time_dict
17

```

```

1 # 去取文章相似度
2 i2i_sim = pickle.load(open(save_path + 'itemcf_i2i_sim.pkl', 'rb'))
3
4 # 相似文章的数量
5 sim_item_topk = 10
6
7 # 召回文章数量
8 recall_item_num = 10
9
10 # 用户热度补全
11 item_topk_click = get_item_topk_click(all_click_df, k=50)
12

```

```

1 # 获取近期点击最多的文章
2 def get_item_topk_click(click_df, k):
3     # value_counts() 相同click_article_id出现次数 并由大到小排序 就是同一篇文章出现次数由
    大到小排列 index取出前50篇
4     topk_click = click_df['click_article_id'].value_counts().index[:k] # 返回的是个列
    表
5     return topk_click
6

```

```

1 for user in tqdm(all_click_df['user_id'].unique()): # 可以理解为去重 就是里面有哪些不同的用户
2     # 返回值是(item,相似系数),....
3     user_recall_items_dict[user] = item_based_recommend(user, user_item_time_dict,
4                                                         i2i_sim,
5                                                         sim_item_topk,
6                                                         recall_item_num, item_topk_click)
5 # user_item_time_dict:{user1: [(item1, time1), (item2, time2)..]...} i2i_sim 相似系数矩阵
6 # sim_item_topk=10 相似文章的数量 recall_item_num = 10 召回文章数量
   item_topk_click 获取近期点击最多的文章

```

给每个用户根据物品的协同过滤推荐文章

```

1 # 基于商品的召回i2i
2 def item_based_recommend(user_id, user_item_time_dict, i2i_sim, sim_item_topk,
3                           recall_item_num, item_topk_click):
4     """
5     基于文章协同过滤的召回
6     :param user_id: 用户id
7     :param user_item_time_dict: 字典, 根据点击时间获取用户的点击文章序列 {user1:
8     [(item1, time1), (item2, time2)..]...}
9     :param i2i_sim: 字典, 文章相似性矩阵
10    :param sim_item_topk: 整数, 选择与当前文章最相似的前k篇文章
11    :param recall_item_num: 整数, 最后的召回文章数量
12    :param item_topk_click: 列表, 点击次数最多的文章列表, 用户召回补全
13    return: 召回的文章列表 {item1:score1, item2: score2...}
14    注意: 基于物品的协同过滤(详细请参考上一期推荐系统基础的组队学习), 在多路召回部分会
15    加上关联规则的召回策略
16    """
17
18    # 获取用户历史交互的文章
19    user_hist_items = user_item_time_dict[user_id] # 其中一个用户浏览信息
20    user_item_time_dict={user1: [(item1, time1), (item2, time2)..]...}
21    # 13036 [(237822, 1508084188196), (257291, 1508084218196)] user_hist_items=
22    [(item1, time1), (item2, time2)..]
23    # [(199198, 1507030404344), (272143, 1507030434344)]
24    user_hist_items_ = {user_id for user_id, _ in user_hist_items} # 应该是item 例如
25    {199198, 272143}
26
27    item_rank = {}
28    for loc, (i, click_time) in enumerate(user_hist_items): # 按照其编号排序 1
29        (item1, time1);2 (item2, time2)
30        # j=item2, item3.. wij=i2i_sim[i][j],i2i_sim[i][j]... 排序后前10篇 注意
31        是排序后的不一定第一个就是user1 只是方便理解
32        for j, wij in sorted(i2i_sim[i].items(), key=lambda x: x[1], reverse=True)
33        [:sim_item_topk]:
34            if j in user_hist_items_: # 如果商品的item{199198, 272143}
35                continue
36

```

```

28         item_rank.setdefault(j, 0) # {item j:0...}
29         item_rank[j] += wij # {item j:i2i_sim[i][j]...}
30
31     # 不足10个, 用热门商品补全
32     if len(item_rank) < recall_item_num:
33         for i, item in enumerate(item_topk_click): # item_topk_click 前50篇
34             if item in item_rank.items(): # 填充的item应该不在原来的列表中 这个意思应
                该是如果用户浏览过, 就pass
35                 continue
36             item_rank[item] = - i - 100 # 随便给个负数就行
37             if len(item_rank) == recall_item_num:
38                 break
39
40     item_rank = sorted(item_rank.items(), key=lambda x: x[1], reverse=True)
41     [:recall_item_num] # 相似度排序 逆序 高到低
42
43     return item_rank

```

enumerate

召回字典转换成df

```

1 # 将字典的形式转换成df
2 user_item_score_list = []
3 # 注释 score是相似系数
4
5 #         item1  item2  item3
6 for user, items in tqdm(user_recall_items_dict.items()): # {user1 : item_rank}
7     # user1  score  score  score
8     for item, score in items: # 构造一个新列表
9         # user2  score  score  score
10        user_item_score_list.append([user, item, score])
11
12 recall_df = pd.DataFrame(user_item_score_list, columns=['user_id',
13 'click_article_id', 'pred_score']) # 我也不懂这里为啥是预测分数。好像全局都没预测

```

Pandas-DataFrame基础知识点

```

1 # 获取测试集
2 tst_click = pd.read_csv(data_path + 'testA_click_log.csv')
3 tst_users = tst_click['user_id'].unique() # 返回不同的user_id
4
5 # 从所有的召回数据中将测试集中的用户选出来
6 tst_recall = recall_df[recall_df['user_id'].isin(tst_users)] # 在user_id列表中查找
   tst_users 返回值=['user_id', 'click_article_id', 'pred_score']
7
8 # 生成提交文件
9 submit(tst_recall, topk=5, model_name='itemcf_baseline')

```

精确查找isin

```

1 # 生成提交文件
2 def submit(recall_df, topk=5, model_name=None):
3     recall_df = recall_df.sort_values(by=['user_id', 'pred_score']) # 排序 如果
   user_id相同按照pre_score排序
4     recall_df['rank'] = recall_df.groupby(['user_id'])
   ['pred_score'].rank(ascending=False, method='first') # 名次 最小的的排名第一
5
6     # 判断是不是每个用户都有5篇文章及以上
7     tmp = recall_df.groupby('user_id').apply(lambda x: x['rank'].max()) # 取最大的名
   次是否大于5
8     assert tmp.min() >= topk
9
10    del recall_df['pred_score']
11    submit = recall_df[recall_df['rank'] <= topk].set_index(['user_id',
   'rank']).unstack(-1).reset_index()
12
13    submit.columns = [int(col) if isinstance(col, int) else col for col in
   submit.columns.droplevel(0)]
14    # 按照提交格式定义列名
15    submit = submit.rename(columns={'': 'user_id', 1: 'article_1', 2: 'article_2',
   3: 'article_3', 4: 'article_4', 5:
   'article_5'})
16
17
18    save_name = save_path + model_name + '_' + datetime.today().strftime('%m-%d') +
   '.csv'
19    submit.to_csv(save_name, index=False, header=True)

```