

In [1]:

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

# Для работы с матрицами
from scipy.sparse import csr_matrix, coo_matrix

# Детерминированные алгоритмы
from implicit.nearest_neighbours import ItemItemRecommender, CosineRecommender, TFIDFRecommender

# Метрики
from implicit.evaluation import train_test_split
from implicit.evaluation import precision_at_k, mean_average_precision_at_k, AUC_at_k, ndcg

```

In [2]:

```

data = pd.read_csv('data/retail_train.csv')
data.head(2)

```

Out[2]:

	user_id	basket_id	day	item_id	quantity	sales_value	store_id	retail_disc	trans_time
0	2375	26984851472	1	1004906	1	1.39	364	-0.6	1631
1	2375	26984851472	1	1033142	1	0.82	364	0.0	1631

In [3]:

```

test_size_weeks = 3

data_train = data[data['week_no'] < data['week_no'].max() - test_size_weeks]
data_test = data[data['week_no'] >= data['week_no'].max() - test_size_weeks]

```

In [4]:

```

popularity = data.groupby('item_id')['sales_value'].sum().reset_index()
popularity[:3]

```

Out[4]:

	item_id	sales_value
0	25671	20.94
1	26081	0.99
2	26093	1.59

In [5]:

```
pb = data_train.groupby('item_id')['sales_value'].sum().reset_index()
pb.columns=['item_id', 'sales']
r_sum = np.sum(pb['sales'])
pb['sales'] = pb['sales']/r_sum
# pb.head(2)
pb['sales'].sum()
```

Out[5]:

0.9999999999999999

## Задание 1. Weighted Random Recommendation

Напишите код для случайных рекомендаций, в которых вероятность рекомендовать товар прямо пропорциональна логарифму продаж

- Можно сэмплировать товары случайно, но пропорционально какому-либо весу
- Например, прямопропорционально популярности. Вес =  $\log(\text{sales\_sum товара})$

## Бейзлайн

Создадим датафрейм с покупками юзеров на тестовом датасете (последние 3 недели)

In [6]:

```
result = data_test.groupby('user_id')['item_id'].unique().reset_index()
result.columns=['user_id', 'actual']
result.head(2)
```

Out[6]:

	user_id	actual
0	1	[821867, 834484, 856942, 865456, 889248, 90795...
1	3	[835476, 851057, 872021, 878302, 879948, 90963...

In [7]:

```
test_users = result.shape[0]
new_test_users = len(set(data_test['user_id']) - set(data_train['user_id']))

print('В тестовом дата сете {} юзеров'.format(test_users))
print('В тестовом дата сете {} новых юзеров'.format(new_test_users))
```

В тестовом дата сете 2042 юзеров

В тестовом дата сете 0 новых юзеров

In [8]:

```
def weighted_random_recommendation(items,pb, n=5):
    """Случайные рекомендации"""
    pb = np.array(pb)
    items = np.array(items)
    recs = np.random.choice(items, size=n, replace=False, p = pb)

    return recs.tolist()
```

In [9]:

```
%%time

items = data_train.item_id.unique()

result['weighted_random_recommendation'] = result['user_id'].apply(lambda x: weighted_random_recommendation(items, pb, n=5))

result.head(2)
```

Wall time: 3.57 s

Out[9]:

	user_id	actual	weighted_random_recommendation
0	1	[821867, 834484, 856942, 865456, 889248, 90795...]	[6754701, 1120526, 12604177, 12581888, 1030056]
1	3	[835476, 851057, 872021, 878302, 879948, 90963...]	[12172410, 987655, 824989, 1067749, 6704400]

## Задание 2. Улучшение бейзлайнов и ItemItem

- Попробуйте улучшить бейзлайны, считая случайный на топ-5000 товаров
- Попробуйте улучшить разные варианты ItemItemRecommender, выбирая число соседей  $K$ .

## Random top

In [10]:

```
popularity = data_train.groupby('item_id')['quantity'].sum().reset_index()
popularity.rename(columns={'quantity': 'n_sold'}, inplace=True)

popularity.head()
```

Out[10]:

	item_id	n_sold
0	25671	6
1	26081	1
2	26093	1
3	26190	1
4	26355	2

In [11]:

```
top_5000 = popularity.sort_values('n_sold', ascending=False).head(5000)#.item_id.tolist()
top_5000
```

Out[11]:

	item_id	n_sold
55470	6534178	190227964
55430	6533889	15978434
55465	6534166	12439291
55576	6544236	2501949
43620	1404121	1562004
...	...	...
51914	5574045	120
10610	864048	119
7777	838842	119
26071	1004001	119
23122	977033	119

5000 rows × 2 columns

In [12]:

```
data_train_top = top_5000.merge(data_train, on='item_id')
```

In [13]:

```
len(data_train_top['item_id'].unique())
```

Out[13]:

5000

In [14]:

```
len(data_train['item_id'].unique())
```

Out[14]:

86865

In [15]:

```
def random_recommendation_top(items, n=5):
    """Случайные рекомендации"""

    items = np.array(items)
    recs = np.random.choice(items, size=n, replace=False)

    return recs.tolist()
```

In [16]:

```
%%time

items = data_train_top.item_id.unique()

result['random_recommendation_top'] = result['user_id'].apply(lambda x: random_recommendation_top(items, n=5))

result.head(2)
```

Wall time: 390 ms

Out[16]:

	user_id	actual	weighted_random_recommendation	random_recommendation_top
0	1	[821867, 834484, 856942, 865456, 889248, 90795...]	[6754701, 1120526, 12604177, 12581888, 1030056]	[934948, 5586955, 1138467, 1128900, 918846]
1	3	[835476, 851057, 872021, 878302, 879948, 90963...]	[12172410, 987655, 824989, 1067749, 6704400]	[987237, 7410344, 868548, 986947, 910161]

## itemitemRecommender

In [17]:

```
top_5000 = popularity.sort_values('n_sold', ascending=False).head(5000).item_id.tolist()
```

In [18]:

```
data_train.head()
```

Out[18]:

	user_id	basket_id	day	item_id	quantity	sales_value	store_id	retail_disc	trans_time
0	2375	26984851472	1	1004906	1	1.39	364	-0.60	1631
1	2375	26984851472	1	1033142	1	0.82	364	0.00	1631
2	2375	26984851472	1	1036325	1	0.99	364	-0.30	1631
3	2375	26984851472	1	1082185	1	1.21	364	0.00	1631
4	2375	26984851472	1	8160430	1	1.50	364	-0.39	1631

In [19]:

```
# Заведём фиктивный item_id (если юзер покупал товары из топ-5000, то он "купил" такой товар)
data_train.loc[ ~ data_train['item_id'].isin(top_5000), 'item_id'] = 6666
data_train.head(100)
```

C:\Users\voron\AppData\Roaming\Python\Python37\site-packages\pandas\core\indexing.py:1720: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

```
self._setitem_single_column(loc, value, pi)
```

Out[19]:

	user_id	basket_id	day	item_id	quantity	sales_value	store_id	retail_disc	trans_time
0	2375	26984851472	1	1004906	1	1.39	364	-0.60	1631
1	2375	26984851472	1	1033142	1	0.82	364	0.00	1631
2	2375	26984851472	1	1036325	1	0.99	364	-0.30	1631
3	2375	26984851472	1	1082185	1	1.21	364	0.00	1631
4	2375	26984851472	1	8160430	1	1.50	364	-0.39	1631
...	...	...	...	...	...	...	...	...	...
95	1060	26985040735	1	9553288	1	8.49	315	0.00	1251
96	1351	26985052379	1	903230	1	0.99	447	-0.30	1955
97	744	26985165432	1	6666	0	0.00	31582	0.00	1119
98	212	26985205886	1	822346	1	1.25	288	-0.34	1341
99	212	26985205886	1	830887	1	2.29	288	-0.70	1341

100 rows × 12 columns

In [20]:

```
user_item_matrix = pd.pivot_table(data_train,
                                   index='user_id', columns='item_id',
                                   values='quantity',
                                   aggfunc='count',
                                   fill_value=0
                                   )

user_item_matrix[user_item_matrix > 0] = 1 # так как в итоге хотим предсказать
user_item_matrix = user_item_matrix.astype(float) # необходимый тип матрицы для implicit

# переведем в формат sparse matrix
sparse_user_item = csr_matrix(user_item_matrix).tocsr()
```

In [21]:

```
user_item_matrix.shape
```

Out[21]:

```
(2499, 5001)
```

In [22]:

```
user_item_matrix.sum().sum() / (user_item_matrix.shape[0] * user_item_matrix.shape[1]) * 10
```

Out[22]:

```
5.33770796861036
```

In [23]:

```
userids = user_item_matrix.index.values
itemids = user_item_matrix.columns.values

matrix_userids = np.arange(len(userids))
matrix_itemids = np.arange(len(itemids))

id_to_itemid = dict(zip(matrix_itemids, itemids))
id_to_userid = dict(zip(matrix_userids, userids))

itemid_to_id = dict(zip(itemids, matrix_itemids))
userid_to_id = dict(zip(userids, matrix_userids))
```

In [44]:

```
%%time
K = 30
model = ItemItemRecommender(K=K, num_threads=4) # K - кол-во ближайших соседей

model.fit(csr_matrix(user_item_matrix).T.tocsr(), # На вход item-user matrix
          show_progress=True)

recs = model.recommend(userid=userid_to_id[2], # userid - id от 0 до N
                       user_items=csr_matrix(user_item_matrix).tocsr(), # на вход user-item matrix
                       N=5, # кол-во рекомендаций
                       filter_already_liked_items=False,
                       filter_items=None,
                       recalculate_user=True)
```

```
0%|          | 0/5001 [00:00<?, ?it/s]
```

Wall time: 1.17 s

In [45]:

recs

Out[45]:

```
[(0, 78679.0),
 (3408, 72173.0),
 (2149, 60612.0),
 (3587, 53104.0),
 (2308, 52506.0)]
```

In [46]:

recs

Out[46]:

```
[(0, 78679.0),
 (3408, 72173.0),
 (2149, 60612.0),
 (3587, 53104.0),
 (2308, 52506.0)]
```

In [47]:

```
[id_to_itemid[rec[0]] for rec in recs]
```

Out[47]:

```
[6666, 1082185, 981760, 1098066, 995242]
```



In [48]:

```
%%time

result['itemitem'+ f'_{str(K)}'] = result['user_id'].apply(lambda user_id: [id_to_itemid[re
    model.recommend(userid=userid_to_id[user_id],
    user_items=sparse_user_item,      # на вход user-item matr
    N=5,
    filter_already_liked_items=False,
    filter_items=None,
    recalculate_user=True)])

result.head()
```

Wall time: 168 ms

Out[48]:

	user_id	actual	weighted_random_recommendation	random_recommendation_top	itemitem
0	1	[821867,			[666
		834484,			108218
		856942,	[6754701, 1120526, 12604177,	[934948, 5586955, 1138467,	99524
		865456,	12581888, 1030056]	1128900, 918846]	102974
		889248,			84036
1	3	90795...			
		[835476,			[666
		851057,			108218
		872021,	[12172410, 987655, 824989, 1067749,	[987237, 7410344, 868548,	109806
		878302,	6704400]	986947, 910161]	653417
2	6	879948,			82624
		90963...			
		[920308,			[666
		926804,			108218
		946489,	[907260, 946484, 2066533, 939252,	[1005902, 920308, 1046131,	98176
3	7	1006718,	12604131]	1051093, 1043128]	99524
		1017061,			102974
		107...			
		[840386,			[666
		889774,			108218
4	8	898068,	[871573, 1651380, 872856, 6423911,	[1112238, 967461, 1045220,	99524
		909714,	9856740]	6979089, 5584847]	102974
		929067,			82624
		95347...			
		[835098,			[666
		872137,			108218
		910439,	[1099628, 847625, 12488017,	[9935616, 1031316, 1128900,	98176
		924610,	12604177, 1354554]	969205, 5563934]	99524
		992977,			102974
		10412...			

In [51]:

```
from metrics import precision_at_k, recall_at_k

for name_col in result.columns[1:]:
    print(f"{round(result.apply(lambda row: precision_at_k(row[name_col], row['actual']), a
```

```
1.0:actual
0.0003:weighted_random_recommendation
0.0063:random_recommendation_top
0.1923:itemitem_1
0.1368:itemitem_5
0.1509:itemitem_10
0.1532:itemitem_15
0.1495:itemitem_30
```