

حَسْنَةٌ لِرَحْمَةِ رَبِّنَا



بازیابی اطلاعات

پروژه: بررسی الگوریتم و پیاده‌سازی سیستم‌های پشن‌هاده‌نده (توصیه‌گر)

آرمان آزادی

۱۴۰۴ زمستان

فهرست مطالب

۴.....	مقدمه	.۱
۵.....	بخش اول: توصیف و تحلیل داده‌ها	.۲
۷.....	بخش دوم: پیاده‌سازی سیستم توصیه‌گر	.۳
۱۱.....	بخش سوم: ارزیابی و تحلیل نتایج	.۴
۱۳.....	نتیجه‌گیری نهایی	.۵

۱. مقدمه

این پژوهه تلاشی است برای درک عمیق یک مفهوم اصلی سیستم‌های بازیابی اطلاعات: سیستم‌های توصیه‌گر که تجربه کاربری را شخصی‌سازی می‌کنند.

ما با یک چالش اصلی در حوزه بازیابی اطلاعات روبرو هستیم: طراحی یک سیستم توصیه‌گر کتاب برای پیشنهاد دادن کتاب‌های مناسب به کاربران بر اساس سلیقه آن‌ها. هدف اصلی این پژوهه، درک عمیق و همچنین پیاده‌سازی عملی روش‌های فیلترینگ مشارکتی (Collaborative Filtering) بر روی یک مجموعه داده واقعی است. تمام مراحل با زبان پایتون و بدون استفاده از کتابخانه‌های آماده‌ی توصیه‌گر و صرفاً با استفاده از جبر خطی و محاسبات ماتریسی پیاده‌سازی شده است. این گزارش شرح کاملی از روند انجام پژوهه، چالش‌هایی که با آن‌ها روبرو شده و نتایجی است که به دست آورده شده است.

۲. بخش اول: توصیف و تحلیل داده‌ها

اولین قدم شناخت دقیق داده‌های ارائه شده در این تمرین مربوط به یک سیستم کتابخوانی است.

ساختار فایل‌های ورودی

مجموعه داده شامل چندین فایل CSV بود که آن‌ها را با استفاده از کتابخانه pandas بررسی کرده:

:**interactions.csv**: این مهم‌ترین فایل پروژه بوده. شامل ستون‌های book_id، user_id و rating است. نکته‌ای که در همان ابتدا توجه را جلب می‌کند، وجود تعداد زیادی امتیاز صفر است. در سیستم‌های توصیه‌گر، امتیاز صفر معمولاً به معنی "عدم علاقه" نیست، بلکه می‌تواند به معنی تعامل ضمنی (Implicit Feedback) مثل کلیک کردن یا دیدن صفحه کتاب باشد. این موضوع در بخش پیاده‌سازی بسیار تأثیرگذار بود.

:**books.csv**: این فایل شامل اطلاعات متنی کتاب‌ها مثل عنوان، نویسنده و ناشر است. اگرچه در روش‌های فیلترینگ مشارکتی (Collaborative Filtering) که در بخش ۳ استفاده کردیم نیازی به این محتوا نداریم، اما این فایل برای درک اینکه آیتم‌ها چه هستند مفید است.

:**test_users.csv** و **users.csv**: لیست کاربران کل سیستم و کاربرانی که باید برای آن‌ها توصیه تولید کنیم.

:**gold_interactions.csv**: این فایل شامل پاسخ‌های صحیح (Ground Truth) برای کاربران تست بود که برای ارزیابی نهایی از آن استفاده شده است.

در (شکل ۱) میتوان کلاس پیکربندی برای قسمت پیاده‌سازی سیستم توصیه‌گر (section three) و در (شکل ۲) برای قسمت ارزیابی سیستم توصیه‌گر (section four) را دید:

```
1 class Config:
2     TRAIN_INTERACTIONS = 'interactions.csv'
3     TEST_USERS = 'test_users.csv'
4
5     # Columns
6     USER_ID_COL = 'user_id'
7     BOOK_ID_COL = 'book_id'
8     RATING_COL = 'rating'
9
10    TOP_N = 10
```

شکل ۱

```
1 class Config:
2     GOLD_FILE = 'gold_interactions.csv'
3
4     USER_BASED_FILE = 'user_based_results.csv'
5     ITEM_BASED_FILE = 'item_based_results.csv'
6
7     # Evaluation constraints
8     K_PRECISION = [5, 10]
9     NDCG_K = 10
```

شکل ۲

چالش پراکندگی داده‌ها (Sparsity)

یکی از اولین بررسی‌های آماری که انجام شده، محاسبه ابعاد ماتریس تعاملات بود.

تعداد کاربران یکتا: ۲۲۸۵

تعداد کتاب‌های یکتا: ۶۸۸۰

تعداد کل خانه‌های ماتریس: $2285 * 6880 = 15,720,800$

اما تعداد رکوردهای موجود در فایل interactions.csv بسیار کمتر از این مقدار بود. این یعنی ما با یک ماتریس بسیار تنک (Sparse) روبرو هستیم. این پراکندگی باعث می‌شود که بسیاری از الگوریتم‌های ساده نتوانند همسایگان مشابهی پیدا کنند، چون همپوشانی بین کاربران بسیار کم است.

۳. بخش دوم: پیاده‌سازی سیستم توصیه‌گر

هدف این بود که برای کاربرانی که در `test_users.csv` مشخص شده‌اند، ۱۰ کتاب پیشنهاد دهیم. طبق خواسته پروژه، از روش‌های مبتنی بر Collaborative Filtering استفاده شده است.

رویکرد کلی و ماتریس تعاملات

به جای استفاده از کتابخانه‌های آماده، از جبر خطی با `numpy` استفاده شده است. ابتدا یک ماتریس تعاملات R با ابعاد $N\{\text{users}\} \times M\{\text{books}\}$ ساخته می‌شود (سطر همان کاربر و ستون کتابها). نکته بسیار مهمی که در اینجا رعایت شده است، مدیریت داده‌های صفر بود. اگر فقط امتیازات غیرصفر را در نظر می‌گرفتیم، بسیاری از اطلاعات از دست می‌رفت. پس ماتریس را با مقادیر اصلی پر کردیم.

```
30 # 4. Build Interaction Matrix R and Seen Mask
31 # R: Stores the ratings (0-10)
32 # seen_mask: Stores True if ANY interaction happened (even if rating is 0)
33 R = np.zeros((n_users, n_books), dtype=np.float32)
34 seen_mask = np.zeros((n_users, n_books), dtype=bool)
35
36 # Efficient filling
37 u_idxs = [user_to_idx[u] for u in interactions[Config.USER_ID_COL]]
38 b_idxs = [book_to_idx[b] for b in interactions[Config.BOOK_ID_COL]]
39 r_vals = interactions[Config.RATING_COL].values
40
41 R[u_idxs, b_idxs] = r_vals
42 seen_mask[u_idxs, b_idxs] = True # Mark as seen regardless of rating value
43
```

شکل ۳

محاسبه شباهت

شباهت کسینوسی (Cosine Similarity) مدنظر است که فرمول آن به صورت برداری: شباهت بین سطر u و سطر v در ماتریس امتیازات محاسبه شد. این کار را با ضرب ماتریسی نرمال‌شده انجام میدهیم تا سرعت اجرا بالا برود.

$$\text{Sim}(U, V) = \frac{U \cdot V}{\|U\| \times \|V\|}$$

```

1 def compute_cosine_similarity(M, kind='user'):
2     """
3     Computes Cosine Similarity: DotProduct(A, B) / (|A|*|B|)
4     """
5     print(f"Calculating {kind}-based Cosine Similarity:")
6
7     if kind == 'user':
8         # Sim(u, v) = R[u] . R[v] / (|R[u]| * |R[v]|)
9         row_norms = np.sqrt(np.sum(M**2, axis=1, keepdims=True))
10        row_norms[row_norms == 0] = 1.0 # Avoid div/0
11        M_norm = M / row_norms
12        sim_matrix = np.dot(M_norm, M_norm.T)
13
14    elif kind == 'item':
15        # Sim(i, j) = R[:,i] . R[:,j] / (|R[:,i]| * |R[:,j]|)
16        col_norms = np.sqrt(np.sum(M**2, axis=0, keepdims=True))
17        col_norms[col_norms == 0] = 1.0
18        M_norm = M / col_norms
19        sim_matrix = np.dot(M_norm.T, M_norm)
20
21    np.fill_diagonal(sim_matrix, 0)
22    return sim_matrix

```

شکل ۴

پیاده‌سازی روش کاربر-محور (User-Based)

ایده: "کاربران شبیه به من چه کتاب‌هایی خوانده‌اند که من نخوانده‌ام؟"

برای این کار مراحل زیر را طی کردم:

1. **محاسبه شباهت: شباهت کسینوسی (Cosine Similarity)** را بین سطرهای ماتریس (کاربران) محاسبه می‌کنیم. امتیاز پیش‌بینی شده برای کاربر u و کتاب i برابر است با میانگین وزن‌دار امتیازاتی که "همسایگان" کاربر u به کتاب i داده‌اند.

2. پیش‌بینی :امتیاز کتاب‌ها را با ضرب ماتریس شباهت در ماتریس امتیازات محاسبه می‌کنیم. این کار معادل میانگین‌گیری وزن‌دار از نظرات همسایگان است.

$$\mathbf{P} = \mathbf{Sim}_{\text{user}} \times \mathbf{R}$$

پیاده‌سازی روش آیتم-محور (Item-Based)

ایده: " کتاب‌های شبیه به کتاب‌هایی که قبلاً خوانده‌ام کدامند؟"

این بار شباهت را بین ستون‌های ماتریس محاسبه می‌کنیم. یعنی دو کتاب چقدر توسط کاربران یکسان، امتیاز مشابه گرفته‌اند.

چالش اصلی در اینجا این بود که به دلیل تنکی ماتریس، احتمال اینکه دو کتاب دقیقاً توسط کاربران مشترکی خوانده شده باشند کم است، که باعث می‌شود شباهت آیتم-آیتم اغلب صفر شود(شکل ۵).

$$\mathbf{P} = \mathbf{R} \times \mathbf{Sim}_{\text{item}}$$

```
24 def predict(R, sim_matrix, kind='user'):
25     print(f"Predicting ratings using {kind}-based method:")
26     if kind == 'user':
27         pred = np.dot(sim_matrix, R)
28     elif kind == 'item':
29         pred = np.dot(R, sim_matrix)
30     return pred
```

شکل ۵

مدیریت داده‌های دیده شده (Seen Items)

یکی از الزامات پروژه این بود که پیشنهادها باید از بین Unseen Items باشند.

برای این کار یک ماتریس ماسک (seen_mask) ایجاد میکنیم. هر خانه‌ای که کاربر در فایل اصلی تعاملی با آن داشت (حتی با امتیاز ۰)، در این ماسک True شد. هنگام تولید لیست نهایی ۱۰ تایی، امتیاز آیتم‌های دیده شده را به منفی بینهایت (-inf) تغییر میدهیم تا در مرتب‌سازی نهایی ظاهر نشوند(شکل ۶).

```

1 def save_recommendations(seen_mask, pred_matrix, user_to_idx, idx_to_book, target_uids, filename):
2     results = []
3
4     print(f"Extracting top {Config.TOP_N} recommendations:")
5
6     for uid in target_uids:
7         if uid not in user_to_idx:
8             results.append({'user_id': uid, 'recommendations': ""})
9             continue
10
11         u_idx = user_to_idx[uid]
12         user_preds = pred_matrix[u_idx]
13
14         # CORRECT FILTERING:
15         # Use the seen_mask to filter out ALL interacted items (including 0 ratings)
16         user_seen_bool = seen_mask[u_idx]
17         user_preds[user_seen_bool] = -np.inf
18
19         # Get Top N Indices
20         top_indices = np.argsort(user_preds)[-Config.TOP_N:][::-1]
21
22         # Map back to IDs
23         recs = [idx_to_book[i] for i in top_indices]
```

شکل ۶

ذخیره‌سازی نتایج

در نهایت نتایج مربوط به پیشنهادات را در دو فایل جداگانه item_based_results.csv و user_based_results.csv ذخیره میکنم که بتوان در مرحله ارزیابی آن را با مرجع درستی مقایسه کرد.

۴. بخش سوم: ارزیابی و تحلیل نتایج

پس از اجرای مدل‌ها و تولید فایل‌های item_based_results.csv و user_based_results.csv، نوبت به سنجش عملکرد رسید.

معیارهای ارزیابی

اسکریپت evaluation.ipynb معیارهای زیر را محاسبه می‌کند:

- **Precision@K (P@5, P@10)**: دقت سیستم در ۵ و ۱۰ پیشنهاد اول. یعنی چند درصد پیشنهادها واقعاً در لیست gold کاربر بوده‌اند.
- **NDCG**: معیاری که رتبه پیشنهاد را هم در نظر می‌گیرد. اگر کتاب محبوب کاربر را در رتبه ۱ بیاوریم، امتیاز بیشتری دارد تا رتبه ۱۰.
- **Spearman Correlation**: همبستگی رتبه‌بندی ما با رتبه‌بندی واقعی.

نتایج عددی به دست آمده

جدول ۱

روش آیتم-محور (Item-Based)	روش کاربر-محور (User-Based)	معیار (Metric)
0.0056	0.0384	P@5
0.0060	0.0298	P@10
0.0055	0.0384	NDCG
-0.0040	-0.0026	Spearman

```

--- Evaluating User-Based CF ---
Results for User-Based CF:
Mean P@5:      0.0384
Mean P@10:     0.0298
Mean NDCG:     0.0384
Mean Spearman: -0.0026

--- Evaluating Item-Based CF ---
Results for Item-Based CF:
Mean P@5:      0.0056
Mean P@10:     0.0060
Mean NDCG:     0.0055
Mean Spearman: -0.0040

```

شکل ۷

تحلیل چرایی تفاوت عملکرد روش‌ها

همانطور که می‌بینید، روش User-Based حدود ۶ تا ۷ برابر بهتر از Item-Based عمل کرده است. چرا؟

۱. مشکل همپوشانی آیتم‌ها: تعداد کتاب‌ها (۶۸۸۰) نسبت به تعداد کاربران زیاد است. احتمال اینکه دو کتاب خاص دقیقاً توسط یک مجموعه کاربر خوانده شده باشند بسیار پایین است. این باعث می‌شود ماتریس شباهت آیتم-آیتم بسیار خالی باشد و الگوریتم نتواند همسایگان درستی برای کتاب‌ها پیدا کند.
۲. رفتار کاربران: کاربران معمولاً تمایل دارند کتاب‌های محبوب (Popular) را بخوانند. پیدا کردن کاربرانی که سلیقه مشابه در کتاب‌های پرفروش دارند راحت‌تر است. روش User-Based از این "خوشه‌های رفتاری" بهتر استفاده می‌کند.
۳. تأثیر داده‌های صفر: استفاده از داده‌های ضمنی (امتیاز ۰) به User-Based کمک کرد تا پروفایل کاربران را پرتر کند، اما در Item-Based به دلیل تنوع زیاد کتاب‌ها، تأثیر کمتری داشت.

۵. نتیجه‌گیری نهایی

در این پژوهه، ما چرخه کامل یک سیستم توصیه‌گر را طی کردیم.

- در بخش توصیه‌گر، دریافتیم که برای داده‌های پراکنده مثل کتاب، رویکرد User-Based Collaborative در پیش‌پردازش درست داده‌ها (مثل مدیریت امتیازهای صفر و فیلتر کردن کتاب‌های دیده شده) تاثیر حیاتی بر کیفیت خروجی نهایی دارد. اگر این موارد را رعایت نمی‌کردیم، کارایی بسیار بالاتری دارد.
- همچنین یاد گرفتیم که پیش‌پردازش درست داده‌ها (مثل مدیریت امتیازهای صفر و فیلتر کردن کتاب‌های دیده شده) تاثیر حیاتی بر کیفیت خروجی نهایی دارد. اگر این موارد را رعایت نمی‌کردیم، Precision سیستم احتمالاً زیر ۱٪ باقی می‌ماند.

این پژوهه نشان داد که در دنیای واقعی، انتخاب الگوریتم مناسب (کاربر-محور در برابر آیتم-محور) کاملاً وابسته به چگالی و ساختار داده‌هاست و یک راه حل همیشگی وجود ندارد.