

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

**ارائه روش جدید در سامانه‌های توصیه‌گر  
با رویکرد پالایش گروهی**

کارشناسی ارشد مهندسی کامپیوتر - نرم‌افزار

امیرحسین قرائتی

اساتید راهنما

دکتر سید رسول موسوی

دکتر چنگیز اصلاح‌چی

۱۳۹۳



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر گرایش نرم افزار

آقای امیرحسین قرائتی تحت عنوان

**ارائه روش جدید در سامانه‌های توصیه‌گر با رویکرد پالایش گروهی**

در تاریخ ۹۳/۱۰/۲۴ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهایی قرار گرفت.

دکتر سید رسول موسوی

۱-استاد راهنمای اول پایان نامه

دکتر چنگیز اصلاح‌چی

۲-استاد راهنمای دوم پایان نامه

دکتر عبدالرضا میرزائی

۳-استاد داور

دکتر محمدعلی خسروی فرد

سرپرست تحصیلات تکمیلی دانشکده

کلیه حقوق مادی مترتب بر نتایج  
مطالعات، ابتکارات و نوآوری‌های ناشی  
از تحقیق موضوع این پایان‌نامه (رساله)  
متعلق به دانشگاه صنعتی اصفهان است.

## فهرست

|    |                                                                 |
|----|-----------------------------------------------------------------|
| ۱  | چکیده                                                           |
| ۲  | فصل اول: مقدمه‌ای بر سامانه‌های توصیه گر                        |
| ۲  | ۱-۱ معرفی سامانه توصیه گر                                       |
| ۳  | ۲-۱ دلایل اهمیت سامانه توصیه گر                                 |
| ۴  | ۳-۱ کاربردهای سامانه توصیه گر                                   |
| ۴  | ۴-۱ انواع سامانه‌های توصیه گر                                   |
| ۴  | ۱-۴-۱ رویکرد پالایش گروهی                                       |
| ۵  | ۲-۴-۱ مشکلات رویکرد پالایش گروهی                                |
| ۶  | ۳-۴-۱ رویکرد محتوا محور                                         |
| ۷  | ۴-۴-۱ رویکرد ترکیبی                                             |
| ۸  | ۵-۴-۱ مقایسه رویکردهای پالایش گروهی و محتوا محور                |
| ۹  | ۵-۱ اهداف این پایان‌نامه                                        |
| ۱۰ | ۶-۱ نوآوری‌ها در این پایان‌نامه                                 |
| ۱۰ | ۷-۱ مروری بر فصل‌های این پایان‌نامه                             |
| ۱۲ | فصل دوم: بررسی روش‌هایی که با رویکرد پالایش گروهی ارائه شده‌اند |
| ۱۲ | ۱-۲ مقدمه                                                       |
| ۱۳ | ۲-۲ الگوریتم‌های کاربر محور و آیتم محور                         |
| ۱۴ | ۳-۲ بررسی چند روش پالایش گروهی بر اساس حافظه                    |
| ۱۴ | ۴-۲ الگوریتم‌های کاربر محور بر اساس همبستگی پیرسون              |
| ۱۶ | ۵-۲ الگوریتم‌های استنتاج بر مبنای گراف دوبخشی                   |
| ۱۷ | ۱-۵-۲ تصویرسازی گراف دوبخشی                                     |

|    |                                                                 |
|----|-----------------------------------------------------------------|
| ۲۲ | ۲-۵-۲ تصویرسازی گراف دوبخشی و ماتریس وزن                        |
| ۲۶ | ۲-۵-۳ چند روش بر مبنای گراف دوبخشی                              |
| ۲۷ | ۲-۵-۴ پالایش گروهی استاندارد                                    |
| ۲۷ | ۲-۵-۵ استنتاج بر مبنای شبکه NBI                                 |
| ۳۲ | ۲-۵-۶ استنتاج براساس گراف دوبخشی وزن دار یا NBIw                |
| ۳۵ | ۲-۵-۷ استنتاج بهبود یافته براساس شبکه INBI                      |
| ۳۸ | ۲-۶ بررسی پیچیدگی زمانی روش های کاربر محور و آیتم محور          |
| ۴۰ | فصل سوم: روش جدید در سامانه های توصیه گر با رویکرد پالایش گروهی |
| ۴۰ | ۳-۱ مقدمه                                                       |
| ۴۱ | ۳-۲ ماتریس ارزش                                                 |
| ۴۲ | ۳-۳ ضریب تشابه کاربر تصادفی با کاربر فعال                       |
| ۴۴ | ۳-۴ ضریب محبوبیت یک آیتم در همسایگی کاربر فعال                  |
| ۴۵ | ۳-۵ محاسبه عناصر ماتریس ارزش                                    |
| ۵۰ | فصل چهارم: ارزیابی تجربی                                        |
| ۵۰ | ۴-۱ مقدمه                                                       |
| ۵۱ | ۴-۲ مجموعه داده ها                                              |
| ۵۳ | ۴-۳ معیار ارزیابی                                               |
| ۵۴ | ۴-۴ تعریف معیار دقت جدید                                        |
| ۵۵ | ۴-۵ نتایج ارزیابی                                               |
| ۵۸ | فصل پنجم: نتیجه گیری                                            |
| ۵۸ | ۵-۱ جمع بندی                                                    |
| ۵۹ | ۵-۲ پیشنهاد برای بهبود روش جدید                                 |

۳-۵ زمینه فعالیت برای تحقیقات آینده ..... ۶۰

۴-۵ چالش‌های موجود در سامانه‌های توصیه‌گر ..... ۶۱

مراجع ..... ۶۲

## چکیده

سامانه‌های توصیه‌گر، ابزارهای نرم افزاری و تکنیک‌هایی هستند که برای کاربران پیشنهادهایی از آیت‌ها تهیه می‌کنند. امروزه با رشد سریع وب و اطلاعات موجود بر روی آن، سامانه‌های توصیه‌گر توجه زیادی به خود جلب کرده‌اند. بسیاری از وب‌سایت‌های معتبر برای ارائه خدمات به کاربران خود از یک توصیه‌گر بهره می‌برند. به‌طور کلی، یک سامانه توصیه‌گر با رویکرد پالایش گروهی، از کاربران و آیت‌ها تشکیل شده است، به شکلی که هر کاربر تعدادی از آیت‌ها را در کتابخانه خود جمع‌آوری کرده و سامانه توصیه‌گر با توجه به این آیت‌ها و آیت‌هایی که دیگر کاربران جمع‌آوری کرده‌اند، آیت‌های جدیدی که کاربر تا به حال جمع‌آوری نکرده و احتمالاً به آن‌ها علاقمند باشد، به او پیشنهاد می‌دهد.

تاکنون برای توسعه توصیه‌گرها از تکنیک‌ها و تکنولوژی‌های مختلفی استفاده شده است. ما در این تحقیق روش جدیدی در سامانه‌های توصیه‌گر با رویکرد پالایش گروهی ارائه می‌دهیم که بر اساس حافظه است و از تشابه بین دو کاربر برای تولید پیشنهاد استفاده می‌کند، همچنین با تعریف یک همسایگی وزن‌دار در بین آیت‌ها، یک پس‌زمینه کلی از علاقه به آیت‌ها را در نظر می‌گیریم که باعث افزایش دقت پیشگویی الگوریتم جدید می‌شود.

مجموعه داده ویژه‌ای با نام MovieLens برای ارزیابی توصیه‌گرها ایجاد شده، که اکثر توصیه‌گرها را با این مجموعه داده ارزیابی کرده‌اند. ما روش پیشنهادی در این پایان‌نامه را با روشی کاربر محور بر اساس همبستگی پیرسون که یکی از پرکاربردترین روش‌ها در بین روش‌های کاربر محور است و روش‌های NBI، NBIw و INBI، که بر اساس دانش فعلی ما نتایج ارزیابی قابل قبولی در بین دیگر روش‌های آیت محور دارند، مقایسه کرده‌ایم. نتایج آزمایشات با توجه به معیارهای ارزیابی، نشان می‌دهد روش پیشنهادی نتایج ارزیابی بهتری، نسبت به این روش‌ها دارد.



## فصل اول

### مقدمه‌ای بر سامانه‌های توصیه‌گر

#### ۱-۱ معرفی سامانه توصیه‌گر

سامانه توصیه‌گر<sup>۱</sup> (RS) یا سامانه پیشنهادگر، ابزارهای نرم‌افزاری و تکنیک‌هایی است که برای کاربر پیشنهادهایی از آیتم‌ها برای استفاده تهیه می‌کند. پیشنهادها مربوط به فرآیندهای مختلف تصمیم‌گیری هستند، مانند خریدن کدام کالا، گوش دادن به کدام موسیقی، یا خواندن کدام اخبار آنلاین. به‌طور کلی، از واژه آیتم برای معنی کردن آنچه سیستم توصیه می‌کند، استفاده می‌شود.

سامانه توصیه‌گر به‌طور مثال برای افرادی که تجربه شخصی و قدرت مقایسه کافی، برای انتخاب بین تعداد فراوان آیتم‌های موجود در یک وب‌سایت را ندارند، مفید است. برای مثال، سامانه توصیه‌گر در یک وب‌سایت فروش کتاب به کاربران کمک می‌کند کدام کتاب را برای خواندن انتخاب کنند. امروزه در سایت‌های پرتعداد مانند Amazon.com از یک سامانه توصیه‌گر استفاده می‌کنند تا سایت را برای هر کاربر شخصی‌سازی<sup>۲</sup> کنند. سامانه‌های توصیه‌گر شخصی‌سازی نشده هم وجود دارند، مانند لیست پرفروش‌ترین کتاب‌ها. این توصیه‌گرها توسعه به مراتب

---

۱- Recommendation system

۲- Personalized

ساده‌تری دارند. با این‌که، این‌گونه توصیه‌گرها ممکن است در بعضی شرایط مفید باشند، ولی توصیه‌گرهای شخصی‌سازی نشده موضوع اصلی تحقیقات سامانه‌های توصیه‌گر نیستند.

در ساده‌ترین حالت، توصیه‌های شخصی‌سازی شده به شکل یک لیست امتیاز<sup>۱</sup> از آیتم‌ها هستند. در این لیست امتیاز سامانه توصیه‌گر تلاش می‌کند، مناسب‌ترین محصول یا سرویس را برای کاربر بر اساس ترجیحات او، پیشگویی کند. برای انجام چنین کاری، سامانه توصیه‌گر باید ترجیحات کاربر را که به طور معمول از امتیازی که کاربر به یک آیتم می‌دهد برداشت می‌شود، جمع‌آوری کند [۱].

توسعه سامانه‌های توصیه‌گر از یک ایده ساده شروع شد: اغلب افراد تصمیمات روزانه خود را بر اساس توصیه‌هایی که دیگران به آن‌ها می‌دهند، می‌گیرند. این توصیه‌ها می‌تواند به صورت کلامی باشد، و یا به شکل توصیه‌نامه، یا مقالات بررسی کتاب و فیلم که در روزنامه‌ها چاپ می‌شود. سامانه‌های توصیه‌گر سعی می‌کنند به این فرآیند طبیعی اجتماعی کمک کنند و آن را تقویت کنند [۲].

## ۱-۲ دلایل اهمیت سامانه‌های توصیه‌گر

در این سال‌ها علاقه به سامانه‌های توصیه‌گر افزایش چشم‌گیری داشته است، از جمله می‌توان به مواردی که در ادامه آورده می‌شود اشاره کرد.

رشد انفجاری اطلاعات موجود در وب و رشد سریع سرویس‌های تجارت الکترونیکی<sup>۲</sup> باعث ابهام کاربر در انتخاب بین گزینه‌های موجود می‌شود. سامانه توصیه‌گر نقش مهمی در سایت‌های معتبر اینترنتی از جمله Amazon.com، Netflix، YouTube و ... بازی می‌کند. برخی شرکت‌های رسانه‌ای سرمایه‌گذاری‌های کلانی بر روی توسعه سامانه‌های توصیه‌گر به‌عنوان بخشی از سرویسی که به مشتریان ارائه می‌دهند، می‌کنند. در حقیقت، افزایش فروش آیتم‌ها، دلیل اصلی استفاده از این تکنولوژی توسط ارائه‌دهنده‌های سرویس‌های مختلف بر روی اینترنت می‌باشد [۳].

از طرف دیگر، موارد زیادی تحقیق و توسعه در زمینه‌ی سامانه توصیه‌گر انجام می‌شود که در ژورنال‌های آکادمیک به چاپ می‌رسند؛ و همچنین برگزاری کنفرانس‌ها و کارگاه‌های آموزشی مربوط به این زمینه، بیانگر اهمیت موضوع سامانه توصیه‌گر بین محققان در این روزها می‌باشد [۱].

---

۱- Ranked list

۲- E-commerce

### ۳-۱ کاربردهای سامانه توصیه‌گر

در [۴] محققین یک رده‌بندی<sup>۱</sup> از سامانه‌های توصیه‌گر موجود در اینترنت ارائه دادند، سامانه‌های توصیه‌گر را از ابعاد گوناگون بررسی کرده و آن‌ها را دسته‌بندی کرده‌اند. با توجه به این رده‌بندی و بررسی سامانه‌های توصیه‌گر موجود بر روی اینترنت، می‌توان کاربردهای زیر را برای آن‌ها نام برد:

- تجارت الکترونیک: به کاربران آیتم‌هایی برای خرید بر روی یک فروشگاه آنلاین پیشنهاد می‌دهند مانند Amazon.com.
- سرگرمی: پیشنهاد فیلم و موسیقی به کاربران مطابق سلیقه آن‌ها مانند Netflix و IMDB.
- شبکه‌های اجتماعی: پیشنهاد افراد جدید برای ارتباط، مبتنی بر دوستان مشترک مانند Facebook و....

### ۴-۱ انواع سامانه‌های توصیه‌گر

سامانه‌های توصیه‌گر به‌طور کلی به سه دسته تقسیم می‌شوند؛ که شامل: ۱-پالایش گروهی<sup>۲</sup> (CF)، ۲-محتوا محور<sup>۳</sup>، و گونه سومی تحت عنوان سامانه توصیه‌گر ترکیبی<sup>۴</sup> است.

#### ۴-۱-۱ رویکرد پالایش گروهی

استفاده از الگوریتم‌های پالایش گروهی، یک رویکرد به سامانه‌های توصیه‌گر است که استفاده‌های بسیار زیادی دارد. رویکرد پالایش گروهی براساس جمع‌آوری و بررسی مقدار زیادی از اطلاعات رفتارها، فعالیت‌ها یا ترجیحات کاربران است، تا با توجه به شباهت کاربران، آنچه کاربران انتخاب خواهند کرد را پیشگویی کند. از فواید اصلی رویکرد پالایش گروهی به این نکته می‌توان اشاره کرد که نیاز به محتویات پیچیده و قابل فهم برای ماشین نیست، در نتیجه این قابلیت را دارد تا آیتم‌های پیچیده مانند فیلم‌ها را، بدون احتیاج به فهمیدن خود آیتم، به شکل دقیقی پیشنهاد دهد [۵].

در این رویکرد فرض بر این است که اگر کاربر A دارای سلیقه‌ی مشابهی با کاربر B باشد، در مقایسه با یک کاربر تصادفی، در یک موضوع مشترک X، دارای عقیده مشابهی با B خواهد بود. به‌طور کلی، پالایش گروهی بر اساس این

---

۱- Taxonomy  
 ۲- Collaborative Filtering  
 ۳- Content Base  
 ۴- Hybrid

فرض است که افرادی که در گذشته با هم هم‌نظر بوده‌اند در آینده هم خواهند بود و اگر در گذشته آیت‌های مشابهی را دوست داشته‌اند، به احتمال زیاد در آینده نیز آیت‌های مشابهی را دوست خواهند داشت [۶].

در رویکرد پالایش گروهی جمع‌آوری اطلاعات به دو شکل صریح<sup>۱</sup> یا ضمنی<sup>۲</sup> انجام می‌شود. در شکل صریح کاربر تمایلات خود را در مورد یک آیت به طور مستقیم بیان می‌کند، معمولاً به شکل امتیاز عددی بین ۱ تا ۵. در شکل ضمنی رفتار و یا انتخاب‌های کاربر برای بدست آوردن تمایلات او ترجمه می‌شوند، برای مثال سابقه خرید کاربر در یک سایت فروش آنلاین.

این رویکرد یک مدل از تمایلات کاربر فعال<sup>۳</sup> و تمایلات دیگر کاربران می‌سازد (با اطلاعاتی که به شکل صریح یا ضمنی جمع‌آوری می‌کند). سپس، از مقایسه مدل‌ها، آیت‌هایی که کاربر ممکن است به آنها علاقه داشته باشد (یا امتیازی که به آنها می‌دهد) را پیشگویی می‌کند. در این جا منظور از کاربر فعال، کاربری است که می‌خواهیم به او پیشنهاد بدهیم.

Amazon.com، Facebook، LinkedIn و ... مثال‌های تجاری و غیر تجاری هستند که از سامانه توصیه‌گری با رویکرد پالایش گروهی استفاده می‌کنند.

### ۱-۴-۲ مشکلات رویکرد پالایش گروهی

سامانه‌های توصیه‌گر با رویکرد پالایش گروهی اغلب با مشکلاتی همچون شروع سرد<sup>۴</sup>، مقیاس‌پذیری<sup>۵</sup> و پراکندگی<sup>۶</sup> مواجه هستند.

مشکل شروع سرد وقتی پیش می‌آید که سامانه توصیه‌گر بخواهد برای کاربری که تاکنون هیچ آیتی را انتخاب نکرده است، پیشنهاد تولید کند. در چنین شرایطی از آنجایی که هیچ اطلاعاتی راجع به تمایلات کاربر در اختیار نداریم رویکرد پالایش گروهی به تنهایی از پس این مشکل بر نمی‌آید. در عمل از ترکیب این رویکرد با رویکرد محتوا محور برای حل این مشکل استفاده می‌شود [۷].

---

۱- Explicit  
 ۲- Implicit  
 ۳- Active User  
 ۴- Cold Start  
 ۵- Scalability  
 ۶- Sparsity

از آنجایی که این سیستم‌ها، اکثراً در محیط‌هایی که میلیون‌ها کاربر و محصول وجود دارد تولید پیشنهاد می‌کنند، به قدرت محاسباتی زیادی برای بررسی سوابق کاربران نیاز دارند [۸]. به‌طور مثال، در سایت Twitter میلیون‌ها کاربر وجود دارد، ولی با به‌کارگیری تکنیک خاصی قدرت محاسباتی مورد نیاز سامانه توصیه‌گر را کاهش می‌دهند، تا سیستم بتواند پیشنهادهای آنلاین برای کاربران تولید کند [۹].

مشکل پراکندگی به این مسئله بر می‌گردد که با وجود میلیون‌ها آیتم، یک فرد به تنهایی نمی‌تواند حتی یک درصد از کل آیتم‌ها را بررسی کند [۸]. در بسیاری از موارد یک آیتم تنها توسط یک کاربر انتخاب شده و یا یک کاربر تنها یک آیتم را انتخاب کرده است. به‌طور مثال، در سایت Amazon.com چند میلیون آیتم وجود دارد که از یک کاربر انتظار نمی‌رود بیشتر از ۱۰،۰۰۰ تا از این آیتم‌ها را بررسی کرده باشد. مشکل پراکندگی، پیشنهاد به کاربر را برای توصیه‌گر دشوار می‌کند.

#### ۱-۴-۳ رویکرد محتوا محور

از آنجایی که تولید پیشنهاد با رویکرد پالایش گروهی سوابق رفتاری کاربران را بررسی می‌کند، به قدرت محاسباتی بالایی برای تولید پیشنهاد احتیاج دارد، و یا ممکن است بحث حریم خصوصی کاربران پیش بیاید. ولی رویکرد محتوا محور، بدون نیاز به تطابق تمایلات کاربران و تنها بر اساس اطلاعاتی که خود آیتم به همراه دارد، می‌تواند پیشنهاد تولید کند [۱۰].

در سیستم‌های توصیه‌گر محتوا محور، از واژه‌های کلیدی برای توصیف آیتم‌ها استفاده می‌شود، به‌علاوه، یک پروفایل برای کاربر ساخته می‌شود تا انواع آیتم‌هایی که او دوست دارد در آن ذخیره شود. به بیان دیگر، این الگوریتم‌ها تلاش می‌کنند آیتم‌هایی را پیشنهاد بدهند که به آیتم‌هایی که کاربر در گذشته انتخاب کرده شباهت دارند. در عمل، چندین آیتم کاندید با آیتم‌هایی که کاربر در گذشته امتیاز داده مقایسه می‌شوند و آیتم‌هایی که بهترین هم‌خوانی را داشته باشند پیشنهاد می‌شوند [۱۱].

برای مثال، با دانستن ژانر فیلم‌ها و این که کاربر به چه ژانری از سینما علاقه دارد، می‌توانیم به کاربر فیلم‌های جدیدی پیشنهاد بدهیم. در کل، سامانه توصیه‌گر محتوا محور به روشی گفته می‌شود، که با مقایسه محتویاتی که یک آیتم را توصیف می‌کند، با محتویات مورد علاقه کاربر، پیشنهاد تولید می‌کند. برای توسعه این رویکرد، یا از روش بازایی

اطلاعات استفاده می‌شود و یا از تکنیک‌های یادگیری ماشین از جمله کلاسبند بیزین<sup>۱</sup>، درخت تصمیم<sup>۲</sup> و شبکه‌های عصبی<sup>۳</sup> [۱۲].

Pandora یک مثال بسیار خوب برای سامانه توصیه‌گر با رویکرد محتوا محور است [۱۳]. Pandora یک ایستگاه رادیویی آنلاین است و برای توصیف آهنگ‌ها از ۴۰۰ صفت متفاوت استفاده می‌کند. پایگاه داده‌های Pandora را پروژه ژنوم موسیقی<sup>۴</sup> می‌نامند. روزانه افراد زیادی در این کمپانی به بررسی آهنگ‌های جدید و مقداردهی به این صفات مشغول هستند تا سیستم بتواند برای کاربران پیشنهادهای بهتری تولید کند.

### ۱-۴-۴ رویکرد ترکیبی

گونه سوم توصیه‌گرها، سیستم‌های ترکیبی هستند که از ترکیبی از دو نوع رویکرد گفته شده استفاده می‌کنند. طراحان، این نوع سیستم‌ها را غالباً به دو منظور با هم ترکیب می‌کنند: ۱- افزایش عملکرد سیستم، ۲- کاهش اثر نقاط ضعف رویکردها وقتی به تنهایی به کار گرفته می‌شوند. تحقیقات اخیر نشان داده که ترکیب رویکردهای پالایش گروهی و محتوا محور می‌تواند در بعضی شرایط کارا تر باشد. رویکرد ترکیبی می‌تواند از راه‌های مختلفی پیاده‌سازی شود، می‌توان پیش‌بینی‌های رویکردهای پالایش گروهی و محتوا محور را جداگانه محاسبه کرد، و سپس آن‌ها را با هم ترکیب کرد، می‌توان توانایی‌های رویکرد محتوا محور را به رویکرد پالایش گروهی اضافه کرد و بالعکس، یا می‌توان رویکردها را در یک مدل یکی کرد. در چندین مطالعه مشاهده‌اتی، کارایی رویکرد ترکیبی را با رویکردهای پالایش گروهی و محتوا محور به تنهایی، مقایسه کرده‌اند و نشان داده‌اند که رویکرد ترکیبی می‌تواند پیشنهادها را دقیق‌تری نسبت به دو رویکرد مذکور به تنهایی، ایجاد کند. همچنین این رویکرد می‌تواند برای غلبه بر مشکلاتی همچون شروع سرد و یا پراکندگی استفاده شود [۱۴].

Netflix یک مثال خوب از رویکرد ترکیبی است [۶]. آن‌ها، هم با مقایسه عادت‌های کاربران در مشاهده فیلم، به کاربرها فیلم پیشنهاد می‌دهند، و همین‌طور فیلم‌هایی را پیشنهاد می‌دهند که مشخصه‌های مشابه با فیلم‌هایی که کاربر به آن‌ها امتیاز بالا داده را دارا است.

۱- Bayesian Classifier

۲- Decision Tree

۳- Artificial Neural Network

۴- Music Genome Project

## ۱-۴-۵ مقایسه رویکردهای پالایش گروهی و محتوا محور

اولین تلاش‌ها برای پالایش اطلاعات، بر اساس محتوا بود. این سیستم‌ها، آیت‌ها را بر اساس محتوای آن‌ها برای پیشنهاد انتخاب می‌کنند. بنابراین، پروفایل کاربر بیان‌کننده محتوایی است که کاربر به آن‌ها علاقه دارد. این گونه از پالایش اطلاعات وقتی کارا است که داده‌های متنی را بازیابی می‌کنیم، به طوری که هر سند<sup>۱</sup> از مجموعه‌ای از کلمات کلیدی تشکیل شده باشد. به هر حال این سیستم‌ها محدودیت‌هایی دارند [۱۵].

اول از همه این که آیت‌ها باید توسط ماشین قابل بررسی باشند. وقتی اطلاعات چند رسانه‌ای را بازیابی می‌کنیم، در شرایطی که ادراک ماشین از محتوا نسبت به ادراک کاربر از محتوا خیلی ضعیف‌تر است، این کار دشوار می‌شود. اگرچه انتساب صفات توسط یک فرد این مشکل را حداقل در این زمینه حل می‌کند، ولی رویکرد محتوا محور در برخورد با بسیاری از اطلاعات امروزی نامناسب است [۱۵].

مشکل دیگر پالایش محتوا محور عدم توانایی در ارزیابی کیفیت یک آیت است. برای مثال، این روش نمی‌تواند تمایزی بین یک مقاله خوب و یک مقاله بد قابل شود اگر هر دو مقاله از لغات یکسانی استفاده کرده باشند. در حقیقت، کیفیت یک آیت یک ویژگی خیلی موضوعی است که بستگی به سلیقه، ایده‌ها، فرهنگ و ... از هر شخص دارد و ممکن است برای ماشین دشوار باشد تا متوجه شود [۱۵].

در نهایت، پالایش محتوا محور راهی ندارد تا آیت‌های غیر مترقبه<sup>۲</sup> که برای کاربر جالب هستند را پیدا کند، این‌ها آیت‌هایی هستند که به طور معمول به پروفایل کاربر مربوط نیستند [۱۵].

سیستم‌های پالایش گروهی کمتر به این مشکلات حساس هستند، از آن‌جایی که به جای محتوای آیت‌ها بر اساس سلیقه کاربران است. این سیستم، آیت‌هایی را پیشنهاد می‌کند که از کاربرانی که سلیقه مشابهی دارند، امتیاز بالایی دریافت کرده باشند. در این تکنیک‌ها، آیت‌ها توسط افراد امتیازدهی می‌شوند و سیستم نیازی ندارد که محتوای آیت‌ها را بررسی کند و کیفیت یا ارزیابی موضوعی یک آیت هم در نظر گرفته شده است [۱۵].

---

۱- document

۲- serendipitous

رویکرد پالایش گروهی سه فایده کلیدی برای پالایش اطلاعات<sup>۱</sup> فراهم می‌کند که هیچ یک از آن‌ها توسط رویکرد محتوا محور فراهم نشده است: (۱) پشتیبانی از پالایش آیت‌هایی که پردازش محتوای آن‌ها برای یک فرآیند خودکار ساده نیست. (۲) توانایی پالایش آیت‌ها بر اساس کیفیت<sup>۲</sup> و سلیقه<sup>۳</sup>. (۳) توانایی تولید پیشنهادهای غیر مترقبه [۱۶].

اول از همه، در پالایش گروهی ارتباط، کیفیت و علاقه به یک آیت در اطلاعات، توسط انسان‌ها مشخص شده است. به‌عنوان مثال، می‌توانیم آیت‌هایی را پالایش کنیم که آنالیز آن‌ها برای کامپیوترها دشوار است، مانند فیلم‌ها، ایده‌ها، احساسات، مردم و سیاستمدارن [۱۶].

دوم، پالایش گروهی می‌تواند در ابعادی فراتر از یک محتوای ساده، اندازه‌گیری کند که یک آیت به چه خوبی نیاز یا علاقه کاربر را برآورده می‌کند که باعث بهبود پالایش اطلاعات می‌شود. انسان‌ها قادر هستند آیت‌ها را در ابعادی مانند کیفیت یا سلیقه بررسی کنند که پردازش آن‌ها برای کامپیوترها خیلی سخت است. یک جستجوگر محتوا محور نشریات می‌تواند تمام مقالات راجع به یک اتفاق را بازگرداند، ولی تنها با ترکیب با یک رویکرد پالایش گروهی می‌تواند مقالات مرتبط به آن موضوع که خوب نوشته شده‌اند را، بازگرداند [۱۶].

در نهایت، یک رویکرد پالایش گروهی در بعضی مواقع پیشنهادهای غیر مترقبه تولید می‌کند، آیت‌هایی را پیشنهاد می‌دهد که برای کاربر خیلی ارزشمند هستند، اما شامل محتوایی که کاربر انتظار آن را داشته است، نیستند. ما متوجه شدیم که پیشنهادهای غیر مترقبه معمولاً در دامنه فیلم اتفاق می‌افتد، با سیستم پالایش گروهی فیلم‌هایی را به‌طور دقیق پیشنهاد می‌دهیم که در شرایط دیگری کاربر آن‌ها را بررسی نمی‌کرد [۱۶].

## ۱-۵ اهداف این پایان‌نامه

در ابتدا هدف ما انجام تحقیقی در زمینه بیوانفورماتیک و موضوع پیشگویی تعاملات بین پروتئین‌ها و داروها<sup>۴</sup> بود. وقتی می‌گوییم پروتئین و دارو با هم تعامل دارند یعنی این دو ساختاری دارند که به هم متصل می‌شوند، ولی از آنجایی که ساختار تمام پروتئین‌ها شناخته شده نیست، اکثر تعاملات بین پروتئین‌ها و داروها نیز ناشناخته مانده است. در تحقیقات مربوط به پیشگویی تعاملات بین پروتئین و دارو هدف پیشگویی تعاملات ناشناخته بر اساس تعاملات شناخته شده است تا محدوده جستجو تعاملات جدید را برای محققان محدود کند. در طول تحقیق به مقاله‌ای برخوردیم که در آن برای پیشگویی تعامل بین پروتئین و دارو از یک سامانه توصیه‌گر که از روش NBI بهره می‌برد، استفاده کرده بود [۱۷]. در آن مقاله، تعاملات شناخته شده بین پروتئین‌ها و داروها را همانند ارتباط‌های کاربران و فیلم‌ها، با یک گراف دوبخشی

۱- Information filtering

۲- quality

۳- taste

۴- Drug-target interaction prediction



مدل می‌کرد و به‌عنوان ورودی به یک توصیه‌گر NBI می‌داد و در نهایت پیشنهادهای توصیه‌گر که در واقع پیشگویی تعاملات جدید ناشناخته بین پروتئین و دارو بودند را به‌عنوان خروجی می‌داد. ما به این نتیجه رسیدیم که اگر بتوانیم توصیه‌گری با دقت پیشگویی بهتر از NBI توسعه دهیم، قادر خواهیم بود پیشگویی دقیق‌تری از تعاملات بین داروها و پروتئین‌ها داشته باشیم.

با توجه به ماهیت مسئله پیشگویی تعاملات بین پروتئین‌ها و داروها باید از توصیه‌گری با رویکرد پالایش گروهی استفاده شود، زیرا اطلاعات کاملی راجع به محتوای پروتئین‌ها یا همان ساختار آن‌ها نداریم و تنها اطلاعاتی راجع به برخی تعاملات آن‌ها با داروهای مختلف را داریم. همان‌طور که گفتیم رویکرد پالایش گروهی مشکلات اساسی خود را دارد ولی در کاربردی که مد نظر ما بود این مسایل خیلی مشکل‌ساز نبودند و مسئله اصلی تولید پیشنهاد با دقت بیشتر بود. در نتیجه هدف ما به ارائه روش جدیدی در سامانه‌های توصیه‌گر با رویکرد پالایش گروهی که دقت پیشگویی بهتری از روش‌های موجود شناخته شده داشته باشد، تغییر کرد.

#### ۱-۶ نوآوری‌ها در این پایان‌نامه

در این تحقیق روش جدیدی در سامانه‌های توصیه‌گر با رویکرد پالایش گروهی ارائه شده که در فصل سوم بیشتر توضیح داده خواهد شد. در این روش که بر اساس حافظه است ما از تشابه بین کاربران استفاده می‌کنیم. یکی از نوآوری‌های این تحقیق این است که برای اندازه‌گیری تشابه میان دو کاربر فرمول جدیدی ارائه می‌کنیم که کاستی‌های فرمول‌هایی که در روش‌های قدیمی استفاده می‌شود را پوشش می‌دهد. در نوآوری دیگری نسبت به روش‌های قدیمی بر اساس تشابه بین دو کاربر، با استفاده از یک همسایگی وزن‌دار در بین آیتم‌ها، دقت پیشگویی الگوریتم خود را افزایش می‌دهیم. این همسایگی وزن‌دار بیان‌گر پس‌زمینه‌ای کلی از علاقه به آیتم‌ها است

در ادامه نشان می‌دهیم که معیارهای ارزیابی موجود برای بیان دقیق کارایی الگوریتم‌ها کافی نیست و ما معیار دقت جدیدی تعریف می‌کنیم که بیان دقیق‌تری از کارایی روش‌ها دارد. این معیار دقت جدید بر اساس جایگاه پیشنهادها در لیست پیشنهاد است و به طول لیست پیشنهاد بستگی دارد. این معیار جدید در کنار معیار نرخ برخورد می‌تواند بیان دقیقی از کارایی الگوریتم در یک پیشنهاد با طول مشخص داشته باشد.

#### ۱-۷ مروری بر فصل‌های این پایان‌نامه

در این پایان‌نامه، ابتدا در فصل دوم، روش‌های توصیه‌گرهای موجود با رویکرد پالایش گروهی را بررسی می‌کنیم و انواع آن‌ها را توضیح می‌دهیم. از آنجایی که روش جدید بر اساس حافظه است، تمرکز ما در این فصل نیز بیشتر بر روی روش‌های بر اساس حافظه خواهد بود. از بین این روش‌ها، روش‌هایی که از هبستگی پیرسون برای اندازه‌گیری

تشابه بین دو کاربر استفاده می کنند و کاربر محور هستند را به طور کامل توضیح می دهیم که یکی از پرکاربردترین روش ها از نوع کاربر محور است و سپس روش های استنتاج بر مبنای گراف دوبخشی و سیر تکامل آن ها را مفصل توضیح می دهیم که در بین روش های موجود بر اساس دانش فعلی ما نتایج ارزیابی قابل قبولی را دارند.

پس از آن در فصل سوم، روش جدید را ارائه می دهیم. در این فصل ابتدا معیار جدیدی برای تشابه بین دو کاربر تعریف می کنیم که کاستی های معیارهای پیشین را پوشش می دهد، سپس با تعریف همسایگی وزن دار در بین آیتم ها یک پس زمینه کلی از علاقه به آیتم ها در نظر می گیریم و دقت الگوریتم جدید را در تولید پیشنهاد افزایش می دهیم.

در فصل چهارم نتایج ارزیابی روش جدید را با نتایج روش های توضیح داده شده در فصل دوم مقایسه می کنیم. در این فصل با استفاده از مجموعه داده معروفی که برای این کار ایجاد شده است و معیارهای ارزیابی موجود و معیار دقت جدیدی که تعریف می کنیم، روش جدید را در کنار روش های نام برده ارزیابی می کنیم و نشان می دهیم که روش جدید نتایج ارزیابی بهتری نسبت به روش های نام برده شده دارد.

و در نهایت در فصل پنجم این پایان نامه، ابتدا جمع بندی از آن چه گفتیم می کنیم و سپس پیشنهادهایی برای بهبود روش جدید ارائه می دهیم همچنین کاربردهای الگوریتم جدید در کارهای آتی را بیان می کنیم. در آخر چالش های باقی مانده در زمینه سامانه های توصیه گر را نام می بریم.

## فصل دوم

### بررسی روش‌هایی که با رویکرد پالایش گروهی ارائه شده‌اند

#### ۲-۱ مقدمه

همان‌طور که در فصل پیش گفتیم، پالایش گروهی عادات هر کاربر در استفاده از اطلاعات مختلف را بررسی می‌کند. رویکرد پالایش گروهی فواید بیشتری از رویکرد محتوا محور دارد، مانند سادگی در نصب و توانایی پالایش هرگونه اطلاعات یا کالا بدون در نظر گرفتن محتوای آن‌ها. نتایج تجربی نیز نشان داده‌اند که رویکرد پالایش گروهی در بسیاری از مواقع از رویکرد محتوا محور عملکرد بهتری دارد [۱۸]. در حال حاضر پالایش گروهی موفق‌ترین و پراستفاده‌ترین رویکرد در سامانه‌های توصیه‌گر است [۱۴]. تمرکز ما در این تحقیق نیز بر روی روش‌های رویکرد پالایش گروهی است.

در سیستم‌های بر اساس پالایش گروهی، پروفایل کاربر شامل مجموعه‌ای از امتیازات داده شده به آیتم‌ها است. این امتیازات با پرسش از کاربر می‌توانند به شکل مستقیم دریافت شوند یا به شکل ضمنی با مشاهده تعامل کاربر با سیستم. امتیازدهی معمولاً به شکل یک ارزش یکتا است (آیتم‌های مربوط را نشان می‌دهد) یا یک مقدار دودویی (اجازه می‌دهد بین آیتم‌های خوب و بد تمایز قایل شویم) یا به طور معمول از یک ارزش عددی در یک مقیاس ثابت است [۱۵].

امتیازات کاربر در یک جدول که به نام ماتریس امتیاز شناخته می‌شود ذخیره می‌شود. این جدول برای تولید پیشنهاد پردازش می‌شود. بسته به این که داده‌های ماتریس امتیاز چگونه پردازش شوند دو نوع الگوریتم متفاوت داریم: براساس حافظه و براساس مدل [۱۵].

الگوریتم‌های براساس حافظه از تمام جدول برای محاسبه پیشگویی خود استفاده می‌کنند. به‌طور کلی، آن‌ها از معیارهای تشابه استفاده می‌کنند تا کاربرانی که به کاربر فعال شبیه هستند را انتخاب کنند. سپس، پیشگویی از امتیازاتی که همسایه‌ها داده‌اند محاسبه می‌شود (به همین علت به آن‌ها بر اساس همسایه نیز گفته می‌شود). اکثر این الگوریتم‌ها می‌توانند براساس این که تمرکز بر روی پیدا کردن کاربران مشابه هست یا آیتم‌های مشابه به دو دسته کاربر محور و آیتم محور تقسیم شوند [۱۵].

الگوریتم‌های براساس مدل ابتدا یک مدل می‌سازند که نشان دهنده رفتار کاربر است و سپس امتیازدهی کاربران را پیشگویی می‌کنند. پارامترهای مدل به شکل آفلاین از ماتریس امتیاز برآورده می‌شود. در مقالات رویکردهای متفاوتی برای این روش وجود دارد، بیشتر آن‌ها مربوط به یادگیری ماشین، ماشین بردار پشتیبان (SVM) [۱۹]، بر اساس روش جبر خطی، بررسی فاکتور، خوشه بندی، شبکه‌های عصبی، گراف‌ها یا روش احتمالاتی مانند شبکه‌های بیزین است [۲۰]. به‌طور کلی الگوریتم‌های براساس حافظه ساده‌تر هستند، در حالی که نتایج با دقت قابل قبولی دارند. به هر حال، از آنجایی که باید تمام داده را برای یک پیشنهاد پردازش کنند، آن‌ها مشکلات جدی پراکندگی دارند. با تعداد زیادی از کاربران و آیتم‌ها این الگوریتم‌ها برای سیستم‌های آنلاین که آیتم‌ها را در زمان واقعی پیشنهاد می‌دهند مناسب نیستند. همچنین، آن‌ها بیشتر از روش‌های بر اساس مدل به مشکلات اساسی سامانه‌های توصیه‌گر حساس هستند. این مشکلات توسط الگوریتم‌های بر اساس مدل کاهش یافته است [۱۵].

## ۲-۲ الگوریتم‌های کاربر محور و آیتم محور

الگوریتم‌های کاربر محور، با نام همسایه محور نیز شناخته می‌شوند و یکی از محبوب‌ترین استراتژی‌ها در بین روش‌های پالایش گروهی است. این الگوریتم‌ها از یک فرآیند سه مرحله‌ای تشکیل شده‌اند [۱۵]:

(۱) تشابه بین کاربر فعال و دیگر کاربران را محاسبه می‌کنند.

(۲) یک زیرمجموعه از کاربران (همسایه‌ها) براساس شباهت آن‌ها با کاربر فعال پیدا می‌کنند.