

NWS_RS: شخصی‌سازی پیشنهادات براساس اطلاعات کاربران و معیار شباهت وزن دار جدید

مصطفی خلجی

دانشکده مهندسی کامپیوتر، دانشگاه صنعتی خواجه نصیرالدین طوسی، تهران.

نام نویسنده مسئول:

مصطفی خلجی

چکیده

با توجه به افزایش اطلاعات در فضای اینترنت، سیستم‌های توصیه‌گر به ابزاری توانمند برای پیداکردن سلیقه، راهنمایی و هدایت کاربران به سمت اقلام مورد نیاز، تبدیل شده‌اند. پالایش همکارانه یکی از رویکردهای اصلی این سیستم‌ها در امر پیداکردن کاربران هم سلیقه و پیش‌بینی میزان علائق کاربران به اقلام خاص است. با افزایش تعداد کاربران و اقلام در سیستم، این رویکرد دچار مشکل مقیاس‌پذیری می‌شود. از طرفی با توجه به تنگ بودن ماتریس امتیازات، عملکرد سیستم کاهش می‌یابد. از این رو، هدف این مقاله ارائه یک سیستم توصیه‌گر برای شخصی‌سازی پیشنهادات براساس اطلاعات کاربران به همراه معیارشباهت وزن دار جدید با نام NWS است. با بکارگیری اطلاعات جمعیتی کاربران از قبیل بازه‌ی سنی، می‌توان مشکل مقیاس‌پذیری را مدیریت کرد و پیشنهادات را برای کاربران شخصی‌سازی کرد. سیستم توصیه‌گر پیشنهادی مبتنی بر پالایش همکارانه کاربر محور بوده و با استفاده از درجه اطمینان تشابه کاربران به عنوان یک وزن در روند پیش‌بینی، عملکرد سیستم را ارتقا می‌دهد. نتایج آزمایشات سیستم توصیه‌گر پیشنهادی بر روی مجموعه داده MovieLens صورت گرفته و ارزیابی سیستم با استفاده از معیارهای MAE، Accuracy، Precision، Recall و F1، بیانگر بهبود خطای پیش‌بینی، افزایش عملکرد و کارایی سیستم توصیه‌گر پیشنهادی نسبت به سایر روش‌های پالایش همکارانه‌ای است که از معیارهای شباهت مختلفی در امر پیشنهاددهی، استفاده می‌کنند. حداکثر نرخ خطای بهبود یافته سیستم ۱۷،۴ درصد و ۲۰،۲ درصد به ترتیب برای کاربران ۲۰ تا ۳۹ سال و کاربران ۴۰ تا ۶۰ سال می‌باشد.

واژگان کلیدی: سیستم توصیه‌گر، پالایش همکارانه، معیارشباهت وزن دار جدید NWS، شخصی‌سازی پیشنهادات.

مقدمه

با افزایش اطلاعات در اینترنت، فضای مجازی و خریدهای اینترنتی و تعاملات کاربران با یکدیگر، سیستم‌های توصیه‌گر (RS) جهت هدایت کاربران به سمت سلاقی یا نیازهایی که دارند، در بیست سال اخیر و به ویژه در دهه‌ی اول قرن بیست و یکم مورد مطالعه قرار گرفته‌اند و پژوهش‌های بسیاری در این زمینه انجام شده است [۱]. این سیستم‌ها با توانایی کشف میزان ترجیحات کاربران و پیش‌بینی اولویت‌های آن‌ها، اقلامی را که احتمال می‌رود مورد پسند کاربر باشد را از بین حجم عظیمی از داده‌ها، پالایش کرده و با توصیه آن‌ها، از اتلاف وقت او جلوگیری می‌کند. امروزه سیستم‌های توصیه‌گر در کاربردهای متفاوتی از قبیل تجارت الکترونیک، توصیه موسیقی، فیلم، مقاله و غیره مورد استفاده قرار می‌گیرند [۲].

پالایش همکارانه (CF) یکی از مهم‌ترین رویکردهای بکاررفته در این سیستم‌ها می‌باشد. عملکرد این روش به این طریق است که براساس شباهت بین کاربران یا اقلام، پیش‌بینی و پیشنهاددهی را برای کاربران فعال^۲ انجام می‌دهد. این الگوریتم با استفاده از ماتریس امتیازات کاربر-قلم، کاربران یا اقلام مشابه را برای پیش‌بینی میزان ترجیحات کاربران فعال نسبت به اقلامی که خریداری یا مشاهده نکرده‌اند، پیدا می‌کند و در آخر لیستی از پیشنهادات را به آنها توصیه می‌کند. ماتریس امتیازات کاربر-قلم، متشکل از امتیازات کاربران به هر یک از اقلام است که به طور معمول برخی از کاربران تمایلی به امتیاز دادن به اقلامی که مشاهده یا خریداری کرده‌اند، ندارند. این رویکرد از مشکلاتی همچون تنگی داده^۴ و مقیاس‌پذیری^۵ رنج می‌برند. هنگامی که سیستم توصیه‌گر با کمبود امتیازات از سوی سوابق کاربر مواجه می‌شود، مشکل تنگی داده بوجود می‌آید که تاثیر بسزایی بر افزایش خطای پیش‌بینی و کاهش عملکرد سیستم دارد. از طرفی دیگر با افزایش تعداد کاربران و اقلام، همچنین محاسبه معیار شباهت برای تک‌تک کاربران فعال با کاربران دیگر، سیستم دچار مشکل مقیاس‌پذیری می‌شود [۳]. برای پیدا کردن تشابه کاربران نسبت به کاربران فعال، معیارهای شباهت مختلفی استفاده می‌شوند که عبارتند از: Pearson, Jaccard, SPC, CPC, MSD [۴-۸]. از این رو هدف این مقاله، طراحی یک سیستم توصیه‌گر مبتنی بر رویکرد پالایش همکارانه، برای شخصی‌سازی پیشنهادات براساس اطلاعات جمعیتی کاربران (سن) به همراه معیار شباهت وزن‌دار جدید (NWS) است که بتوان خطای پیش‌بینی و عملکرد سیستم را بهبود داد. باتوجه به افزایش روز افزون کاربران و اقلام در فضای وب، می‌توان با تفکیک کردن کاربران براساس بازه‌سنی، مساله مقیاس‌پذیری را مدیریت و پیشنهاددهی را برای هر کاربر شخصی‌سازی کرد. هنگامی که تعداد اقلام مشترک بین کاربران فعال و همسایه زیاد باشد، کاربر همسایه‌ای با بیشترین اقلام مشترک از نظر سلیقه احتمال دارد بسیار مورد اطمینان باشد. این امر را می‌توان به عنوان یک ضریب وزنی در سیستم پیشنهاددهی استفاده کرد که منجر به بهبود خطای پیش‌بینی، افزایش کارایی و عملکرد سیستم توصیه‌گر پیشنهادی می‌شود.

ساختار مقاله به این صورت می‌باشد که در بخش ۲ به مروری بر کارهای محققین پرداخته می‌شود و در بخش ۳ سیستم توصیه‌گر پیشنهادی معرفی می‌شود. بخش ۴، بخش ارزیابی سیستم توصیه‌گر پیشنهادی می‌باشد که در آن به نتایج آزمایشات و مقایسه با روش‌های دیگر پرداخته می‌شود. در نهایت در بخش آخر نیز نتیجه‌گیری ارائه می‌گردد.

۱- مروری بر کارهای محققین

سیستم‌های توصیه‌گر در ابتدا توسط گلدبرگ و همکارانش معرفی شدند [۹]. سیستم‌های توصیه‌گر ابزاری برای هرچه توانمند کردن کاربران در پیدا کردن اقلام خاص، در فضای مجازی هستند. این سیستم‌ها با پیدا کردن سلاقی کاربران پیشنهاداتی را برایشان مهیا می‌سازند. پالایش همکارانه یکی از پرکاربردترین رویکردهای بکارگرفته شده در پیدا کردن سلاقی کاربران است. این رویکرد با استفاده از معیار شباهت‌های مختلف، کاربرانی که از نظر سلیقه به کاربر فعال نزدیک هستند را پیدا و براساس آن میزان علاقه کاربر فعال به اقلام خاص را پیش‌بینی می‌کند. اگرچه این روش از مشکلاتی از قبیل شروع سرد، تنگی داده و مقیاس‌پذیری رنج می‌برد اما فهم و پیاده‌سازی آنها بسیار راحت است و از مدل‌های پایه در سیستم‌های توصیه‌گر می‌باشند. از این رو روش‌های جدیدی مبتنی بر پالایش همکارانه توسط محققین برای بهبود عملکرد سیستم‌های توصیه‌گر ارائه گردیده است [۱۰].

پالایش همکارانه به نوبه خود دارای دو روش اصلی مدل محور^۷ و حافظه محور^۸ است. رویکرد مدل محور با استفاده از روش‌های هوشمند یادگیری ماشین، رفتار کاربران را نسبت به سوابق خودشان مدل سازی می‌کند. از طرفی دیگر روش حافظه محور، با استفاده از

1. Recommender Systems
2. Collaborative Filtering
3. Active Users
4. Data Sparsity
5. Scalability
6. New Weighted Similarity
7. Model-Based

رویکرد نزدیک همسایه، به دنبال پیدا کردن کاربرانی است که از نظر سلیقه و ترجیحات به کاربر فعال نزدیک باشند [۱]. امروزه بسیاری از محققین با ترکیب این دو روش به حل مشکلات مطرح شده می‌پردازند.

برای بهبود عملکرد سیستم توصیه‌گر، یک معیار شباهت جدید با نام PIP^9 معرفی شد که توانست مشکل شروع سرد را حل کند. در این معیار از اختلاف بین امتیازات دو کاربر به همراه بررسی امتیازات هر دو کاربر در شرایط سازش (توافق^۱) و عدم سازش (عدم توافق) با یکدیگر، استفاده شد. روش مذکور دارای ضریب جریمه بوده و در شرایطی که امتیازات کاربران در وضعیت عدم توافق با یکدیگر باشند، اعمال می‌شود [۱۱]. لیو و همکارانش یک معیار شباهت اکتشافی جدید با نام $NHSM^{11}$ ارائه کردند که بجای اثر شدید^{۱۲} روش [۱۱]، تکنیکی^{۱۳} امتیازات را هنگام انتخاب کاربران همسایه کاربر فعال، در نظر گرفتند [۱۲]. بلوگین و همکارانش روش‌هایی را برای بهبود عملکرد سیستم‌های توصیه‌گر معرفی کردند. آنها از روش‌های وزن‌دار شناخته شده با نام HW^{14} [۱۳] و MW^{15} [۱۴] با عدم بکارگیری از معیارهای شباهت، کاربرانی که از نظر سلیقه به کاربر فعال نزدیک بودند را پیدا و در روند پیش‌بینی سیستم پیشنهادی خود استفاده کردند. سیستم آنها از نظر پیچیدگی زمانی نسبت به سایر روش‌هایی که از معیارهای شباهت استفاده می‌کردند، عملکرد قابل قبولی را در پی داشت [۱۵]. چوی و همکارانش در سال ۲۰۱۳ یک معیار شباهت جدید برای انتخاب همسایگان برای هر قلم فعال در پالایش همکارانه ارائه کردند که امتیاز کاربر به قلم را براساس معیار خود وزن‌دار می‌کردند [۱۶]. جواری و همکارانش در سال ۲۰۱۴ یک سیستم توصیه‌گر مبتنی بر پالایش همکارانه و تخصیص منابع ارائه کردند. آنها با استفاده از روش تخصیص منابع توانستند درجه اعتماد هر کاربر را براساس میزان شباهت بدست آمده، محاسبه کنند و در نتیجه عملکرد سیستم را نسبت به سایر روش‌های معمول بهبود بخشند [۱۷]. ژانگ و همکارانش در سال ۲۰۱۶ یک روشی برای بهبود توانایی پیدا کردن کاربران مشابه و قابل اطمینان به کاربر فعال را ارائه دادند، هدف آنها ارائه یک سیستم توصیه‌گر موثر مبتنی بر مدل برای حل مشکل تنگی داده بود [۱۸].

پارک و همکارانش در سال ۲۰۱۵ یک الگوریتم پالایش همکارانه سریع با استفاده از نزدیک‌ترین گراف همسایه برای کاهش دادن مشکل پیچیدگی زمانی، ارائه کردند. روش آنها RCF^{16} نام دارد که فرآیند پردازش پیدا کردن نزدیک‌ترین همسایه را نسبت به پالایش همکارانه معمول، معکوس می‌کند [۱۹]. خلجی در سال ۱۳۹۶ یک سیستم توصیه‌گر ترکیبی فیلم با استفاده از شبکه عصبی و تخصیص منابع ارائه کرد. او توانست مشکل مقیاس‌پذیری را با استفاده از شبکه عصبی نگاشت خود سازمان‌ده و روش پیش‌بینی پیوند در شبکه‌های اجتماعی، حل نماید. روش او با جداسازی کاربران نسبت به اطلاعات جمعیتی^{۱۷}شان و کشف سلايق هر یک از کاربران در اقلام و ژانرهای خاص، کارایی سیستم را بهبود بخشید [۲۰].

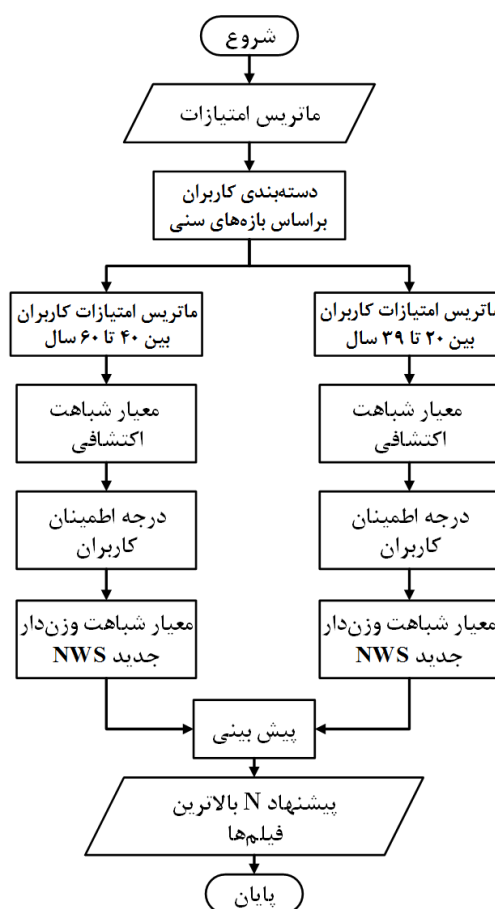
کوهی و همکارش در سال ۲۰۱۶، با بکارگیری خوشه‌بندی فازی و روش فازی‌زدایی Max ، به صورت تخصیص کاربران به تمام خوشه‌ها با درجه عضویت متفاوت و استفاده از معیار شباهت Pearson برای پیدا کردن نزدیک‌ترین همسایه، نشان دادند که عملکرد سیستم آنها نسبت به استفاده از روش‌های K -means و SOM^{18} بهبود بخشیده شده است [۲۱]. بلاسل و همکارانش در سال ۲۰۱۸ یک سیستم توصیه‌گر مقیاس‌پذیر مبتنی بر رویکرد پالایش همکارانه معرفی کردند. آنها با استفاده از الگوریتم خوشه‌بندی جداسازی-ادغام توانستند زمان و دقت سیستم پیشنهادی‌شان را بهبود دهند [۲۲]. کنت و همکارانش در سال ۲۰۱۸ یک روشی برای مشخص کردن انتخاب مرکز اولیه خوشه در عملیات خوشه‌بندی K -means معرفی کردند. روش آنها توانست مشکل تنگی داده را حل کردند [۲۳]. خلجی و دادخواه در سال ۱۳۹۷ یک سیستم توصیه‌گر ترکیبی با نام $FNHSM_HRS$ معرفی کردند. آنها با بکارگیری روش خوشه‌بندی فازی و استفاده از یک معیار شباهت اکتشافی مشکل مقیاس‌پذیری را حل نمودند. سیستم آنها در ابتدا رفتار کاربران را براساس روش‌های فازی مدل کرد، سپس برای پیدا کردن نزدیک‌ترین همسایه از معیار شباهت اکتشافی استفاده کرد [۳]. خداوردی و همکارانش در سال ۱۳۹۷ یک سامانه توصیه‌گر ترکیبی فیلم مبتنی بر خوشه‌بندی و محبوبیت را ارائه کردند. سیستم آنها با استفاده از روش خوشه‌بندی K -means، کاربران مشابه را در خوشه‌های مجزا قرار داد و با بکارگیری روش محبوبیت امتیازات، میزان علاقه کاربر فعال را برای فیلم‌هایی که مشاهده نکرده بود، پیش‌بینی کرد و در آخر روش آنها توانست مشکل مقیاس‌پذیری و تنگی داده را حل کند [۱۰]. خلجی و محمدنژاد در سال ۲۰۱۹ یک

8. Memory-Based
9. Proximity, Impact, Popularity
10. Agreement
11. New Heuristic Similarity Measure
12. Impact
13. Singularity
14. Herlocker Weighting
15. McLaughlin Weighting
16. Reversed CF
17. Demographic Information
18. Self-Organizing Map

سیستم توصیه‌گر ترکیبی فیلم با نام FCNHSMRA_HRS معرفی کردند. سیستم آنها با ترکیب روش‌های مدل محور و حافظه محور در پالایش همکارانه به همراه یکی از روش‌های پیش‌بینی پیوند، توانستند عملکرد سیستم خود را نسبت به روش قبلی خود (FNHSM_HRS) و سایر روش‌های مرسوم بهبود بخشند [۲۴]. وانگ و همکارانش در سال ۲۰۱۹ یک روشی جدید با نام CDIE^{۱۹} معرفی کردند. آنها از روش Co-Clustering برای استخراج همبستگی اقلام و پالایش کردن اختشاش‌ها استفاده کردند. روش آنها توانست مشکل تنگی داده را رفع کند [۲۵].

۲- سیستم توصیه‌گر پیشنهادی

در شکل شماره ۱ ساختار سیستم توصیه‌گر پیشنهادی نشان داده شده است. این ساختار مبتنی بر الگوریتم پالایش همکارانه حافظه محور و روش کاربر محور است.



شکل شماره ۱: ساختار سیستم توصیه‌گر پیشنهادی

سیستم توصیه‌گر پیشنهادی دارای ماتریس امتیازات کاربر-فیلم است که بیانگر تعدادی امتیاز کاربران به فیلم‌ها است. در این مقاله مجموعه کاربران با $U = [u_1, u_2, \dots, u_M]$ ، فیلم‌ها با $I = [i_1, i_2, \dots, i_N]$ و ماتریس امتیازات با نام RMatrix معرفی می‌شوند. سائز RMatrix برابر است با تعداد کاربران \times تعداد فیلم‌ها $(M * N)$.

۱-۲- دسته‌بندی کاربران براساس بازه‌های سنی

در این بخش تمامی کاربران براساس اطلاعات جمعیتی ایشان، دسته‌بندی می‌شوند. اطلاعات جمعیتی کاربران شامل چندین ویژگی از قبیل جنسیت، حرفه(شغل)، میزان تحصیلات، سن و کدپستی است. باتوجه به گذر زمان و افزایش سن، سلیقه افراد تغییرپذیر است. از

این رو، کاربرانی که در بازه سنی جوان دسته‌بندی می‌شوند، از نظر سلیقه و نوع امتیازدهی با کاربرانی که در بازه سنی میانسال قرار دارند، متفاوت است. به طور مثال، کاربران جوان تمایل دارند، فیلم‌هایی با ژانر کنشی مشاهده کنند، اما کاربران میانسال ممکن است به این نوع از ژانر علاقه نشان ندهند و ژانر دیگری را برای خود برگزینند. در سیستم توصیه‌گر پیشنهادی کاربران فقط براساس بازه‌های سنی ۲۰ تا ۳۹ سال و ۴۰ تا ۶۰ سال از همدیگر تفکیک شدند. دسته‌بندی کاربران باعث تسریع در روند پیش‌بینی و پیشنهاددهی می‌شود. علاوه بر این کاربران فعال هنگام ورود به سیستم، تنها می‌توانند از نظرات کاربران هم سن خود برای فیلم‌هایی که مشاهده نموده‌اند، توصیه‌هایی را دریافت کنند. در پایان این بخش، دو ماتریس امتیازات برای بازه‌های سنی ۲۰ تا ۳۹ و ۴۰ تا ۶۰ سال به نام‌های RMatrixA1 و RMatrixA2 با ابعاد $m * n$ از ماتریس RMatrix به وجود می‌آیند. $RMatrixA1 \cup RMatrixA2 = R'Matrix$.

۲-۲- پیدا کردن K نزدیک‌ترین همسایه

سیستم توصیه‌گر پیشنهادی ماتریس‌های RMatrixA1 و RMatrixA2 را به عنوان ورودی دریافت کرده و براساس بازه‌های سنی، ماتریس متعلق به آن کاربرفعال را برای محاسبه K نزدیک‌ترین همسایه انتخاب می‌کند. در گام بعد، سیستم برای پیدا کردن کاربرانی که از نظر سلیقه با کاربرفعال مشابه هستند را با استفاده از معیار شباهت اکتشافی [۱۲] طبق رابطه (۱) محاسبه می‌کند.

$$NHSM_Sim(u, v) = JPSS_Sim(u, v) \cdot URP_Sim(u, v) \quad (1)$$

این معیار دارای دو ضریب اصلی است که هر یک از ضرایب ترکیبی از تعدادی معیارهای شباهت معمول است. از این رو برای محاسبه معیار شباهت اکتشافی مذکور، در ابتدا می‌بایست معیار شباهت $JPSS_Sim(u, v)$ محاسبه شود. این معیار از دو معیار شباهت دیگر مشتق شده است که در روابط (۲) و (۳) ذکر شده است.

$$JPSS_Sim(u, v) = PSS_Sim(u, v) \cdot Jaccard'_Sim(u, v) \quad (2)$$

$$Jaccard'_Sim(u, v) = \frac{|I_u \cap I_v|}{|I_u| \times |I_v|} \quad (3)$$

$|I_u \cap I_v|$ بیانگر تعداد فیلم‌های مشترکی است که کاربر u و v مشاهده نموده‌اند و $|I_u|$ بیانگر تعداد فیلم‌هایی که کاربر فعال (u) امتیاز داده است و $|I_v|$ بیانگر تعداد فیلم‌هایی که کاربر همسایه (v) امتیاز داده است. معیار PSS_Sim از طریق رابطه (۴) بدست می‌آید.

$$PSS_Sim(r_{u,i}, r_{v,i}) = Proximity(r_{u,i}, r_{v,i}) \cdot Significance(r_{u,i}, r_{v,i}) \cdot Singularity(r_{u,i}, r_{v,i}) \quad (4)$$

از طریق روابط (۵) و (۶) و (۷) بدست می‌آید.

$$Proximity(r_{u,i}, r_{v,i}) = 1 - \frac{1}{1 + \exp(-|r_{u,i} - r_{v,i}|)} \quad (5)$$

$$Significance(r_{u,i}, r_{v,i}) = \frac{1}{1 + \exp(-|r_{u,i} - r_{med}| \cdot |r_{v,i} - r_{med}|)} \quad (6)$$

$$Singularity(r_{u,i}, r_{v,i}) = 1 - \frac{1}{1 + \exp(-|\frac{r_{u,i} + r_{v,i}}{2} - \mu_i|)} \quad (7)$$

$r_{u,i}$ امتیاز فیلم i توسط کاربرفعال u و $r_{v,i}$ امتیاز فیلم i توسط کاربر v است. r_{med} هم میانگین امتیازاتی است که کاربر می‌تواند برای امتیازدهی به یک فیلم خاص تخصیص دهد. ماتریس امتیازدهی در سیستم توصیه‌گر پیشنهادی دارای محدوده امتیازات بین اعداد ۱ تا ۵ است که متوسط آن عدد ۳ می‌باشد. μ_i هم میانگین امتیازات فیلم i توسط تمام کاربران است. آخرین مرحله از رابطه (۱)، محاسبه معیار شباهت URP_Sim است که با استفاده از رابطه (۸) بدست می‌آید.

$$URP_Sim(u, v) = 1 - \frac{1}{1 + \exp(-|\mu_u - \mu_v| \cdot |\sigma_u - \sigma_v|)} \quad (8)$$

μ_u میانگین امتیازات و σ_u انحراف معیار کاربر فعال است که مقدار σ_u طبق رابطه (۹) بدست می‌آید.

$$\sigma_u = \sqrt{\sum_{i \in I_u} (r_{u,i} - \bar{r}_u)^2 / |I_u|} \quad (9)$$

اکثراً معیارهای شباهت موجود به طور مستقیم از میزان ترجیحات (امتیازات) کاربران، برای محاسبه تشابه نزدیک‌ترین کاربران همسایه استفاده می‌کنند و در آخر قادر به پیش‌بینی درست تشابه بین کاربران نیستند. از این رو هنگامی که تعداد فیلم‌های مشترک بین کاربران کم باشد، نتایج تشابه بدست آمده کاربران، قابل اطمینان نیست. فرض کنید تعداد فیلم‌های مشترک بین دو کاربر u و v ، ۵ و تشابه آنها نسبت به یکدیگر ۰,۸ باشد و در مقابل تعداد فیلم‌های مشترک بین دو کاربر u و v ، ۱۰۰ و تشابه آنها ۰,۶ باشد، نمی‌توان به اینگونه معیارها برای پیدا کردن K نزدیکترین همسایه اطمینان کرد. بنابراین در این مقاله، فرض می‌شود کاربرانی که تعداد فیلم‌های مشترک زیادی با کاربر فعال دارند، نسبت به کاربرانی که تعداد فیلم‌های مشترک کمتری دارند، از نظر سلیقه به کاربر فعال نزدیک‌تر و قابل اطمینان‌تر هستند و می‌توان به عنوان یک وزن در معیارهای شباهت استفاده کرد. توابع مختلفی برای تبدیل درجه همپوشانی فیلم‌ها به یک وزن معرفی شده‌اند [۱۳] و [۱۴]. در این توابع از یک حد‌آستانه قابل اطمینان برای همپوشانی فیلم‌های مشترک بین کاربران استفاده می‌شود. برای محاسبه درجه اطمینان تشابه کاربران نسبت به یکدیگر از رابطه (۱۰) استفاده می‌شود.

$$Reliability_Degree(u, v) = \frac{\max(\beta, |I_u \cap I_v|)}{\beta} \quad (10)$$

β حد‌آستانه میزان همپوشانی فیلم‌های مشترک بین دو کاربر است. عملکرد این تابع به این صورت است که کاربرانی که تعداد فیلم‌های مشترک بیشتر از حد آستانه داشته باشند دارای اهمیت و وزن بیشتری در روند پیش‌بینی میزان ترجیحات کاربران فعال، هستند. در گام آخر این بخش، با ترکیب معیار شباهت اکتشافی به همراه تابع وزن‌دار (درجه اطمینان)، یک معیار شباهت وزن‌دار جدید به نام NWS طبق رابطه (۱۱) بدست می‌آید.

$$NWS(u, v) = NHSM_Sim(u, v) \cdot Reliability_Degree(u, v) \quad (11)$$

پس از محاسبه معیار شباهت وزن‌دار جدید NWS ، دو ماتریس جداگانه به نام‌های NWS_A1 و NWS_A2 در ابعاد $m * n$ به ترتیب برای کاربران در بازه‌های سنی ۲۰ تا ۳۹ و ۴۰ تا ۶۰ سال ایجاد می‌شوند. این ماتریس‌ها به صورت متقارن هستند و در هر درایه میزان تشابه وزن‌دار کاربران به کاربر فعال را مشخص می‌کنند. به طور مثال اگر u_1 کاربر فعال، u_2 کاربر همسایه و شباهت وزن‌دار بین کاربر u_1 و u_2 ، ۰,۳۶۵ باشد، با توجه متقارن بودن ماتریس‌ها، شباهت بین کاربر u_2 با کاربر فعال u_1 همان ۰,۳۶۵ است. میزان تشابه بدست آمده $NWS(u, v)$ در بازه (۰ تا ۱) است، چون خروجی هر یک از روابط فوق در بازه ۰ تا ۱ است.

۳-۲- پیش‌بینی

در این مرحله تعداد K نزدیک‌ترین کاربران همسایه براساس بالاترین درجه شباهت وزن‌دار به کاربر فعال در بازه‌های سنی مربوطه، مشخص و انتخاب می‌شوند. مقدار K انتخابی برای همسایه‌های کاربر فعال، ۲۰۰ در نظر گرفته شده است. از این رو با توجه به رابطه (۱۲) عملیات پیش‌بینی امتیازات فیلم‌های مشاهده نشده براساس پالایش همکارانه کاربر محور در سیستم توصیه‌گر پیشنهادی، محاسبه می‌شود.

$$Predict(u, i) = \mu_u + \frac{1}{\sum_{j=1}^K |NWS(u, v_j)|} \cdot \sum_{j=1}^K (r_{v_j, i} - \mu_{v_j}) \cdot NWS(u, v_j) \quad (12)$$

به طوری که u کاربر فعال، v_j کاربر همسایه و i فیلمی می‌باشد که سیستم توصیه‌گر پیشنهادی قرار است میزان ترجیح (امتیاز) کاربر فعال را برای آن پیش‌بینی کند، همچنین K تعداد کل کاربران همسایه، $r_{v_j, i}$ امتیاز کاربر v_j به فیلم i و μ_u میانگین امتیازات کاربر فعال است.

۴-۲- پیشنهاددهی

در این بخش لیستی از پیشنهادات براساس روش Top-N برای کاربران فعال ایجاد می‌شود. مقدار N در این مقاله به ترتیب ۵، ۱۰، ۱۵، ۲۰ و ۳۰ در نظر گرفته شده است.

۳- ارزیابی سیستم توصیه‌گر پیشنهادی

داده‌های ورودی سیستم توصیه‌گر پیشنهادی مجموعه داده فیلم MovieLens شامل ۹۴۳ کاربر و ۱۶۸۲ فیلم با ۱۰۰ هزار امتیاز کاربران به فیلم‌ها، در نظر گرفته شده است [۲۴]. بازه امتیازدهی در این مجموعه داده از اعداد ۱ تا ۵ است که به ترتیب ۱ نشان دهنده مورد پسند نبودن و عدد ۵ نشان دهنده مورد پسند بودن کاربر به یک فیلم خاص می‌باشد. تعداد کاربران در بازه سنی ۲۰ تا ۳۹ سال و ۴۰ تا ۶۰ سال در ماتریس امتیازات RMatrixA1 و RMatrixA2 برای ارزیابی به ترتیب ۵۷۷ و ۲۷۱ کاربر می‌باشد. برای ارزیابی عملکرد سیستم، از روش 5-fold Cross-Validation استفاده شده است که ۸۰ درصد داده‌ها برای آموزش سیستم توصیه‌گر پیشنهادی و ۲۰ درصد داده‌ها برای آزمایش سیستم بکارگرفته می‌شود. روش مذکور شامل ۵ مرحله مستقل از هم می‌باشد که در هر مرحله داده‌های آموزش و آزمایش هر Fold به سیستم تزریق می‌شود [۲۵]. ارزیابی سیستم براساس معیارهای میانگین خطای مطلق (MAE)، دقت^{۲۱}، و بازیابی^{۲۲} طبق جدول (۱) و روابط (۱۳-۱۷) بر روی داده‌های آزمایش محاسبه شده است [۲۶]. مقدار حد آستانه برای محاسبه معیارهای ارزیابی ۴ در نظر گرفته شد، مقدار امتیاز ۴ تا ۵ بیانگر مورد پسند بودن (دوست داشتن) و مقدار ۱ تا ۴ بیانگر مورد پسند نبودن (دوست نداشتن) کاربرفعال برای فیلم است. علاوه بر این مقدار حد آستانه β برای محاسبه درجه اطمینان تشابه کاربران در سیستم توصیه‌گر پیشنهادی ۵۰ در نظر گرفته شده است. این سیستم، با ۵ روش معمول و پرکاربرد در این حوزه همچون SPCC, Jaccard, Pearson, MSD و CPC مورد مقایسه و ارزیابی قرار گرفت. نتایج آزمایشات در شکل‌های (۲) تا (۱۱) بیانگر بهبود خطای پیش‌بینی و عملکرد سیستم نسبت به سایر معیارهای شباهت در هر دو بازه سنی است.

جدول ۱: ماتریس اغتشاش

Actual / Predicted	Negative	Positive
Negative	A	B
Positive	C	D

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{r}_{u,i} - r_{u,i}| \quad (13)$$

$$Accuracy = \frac{A + D}{A + B + C + D} \quad (14)$$

$$Precision = \frac{D}{B + D} \quad (15)$$

$$Recall = \frac{D}{C + D} \quad (16)$$

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (17)$$

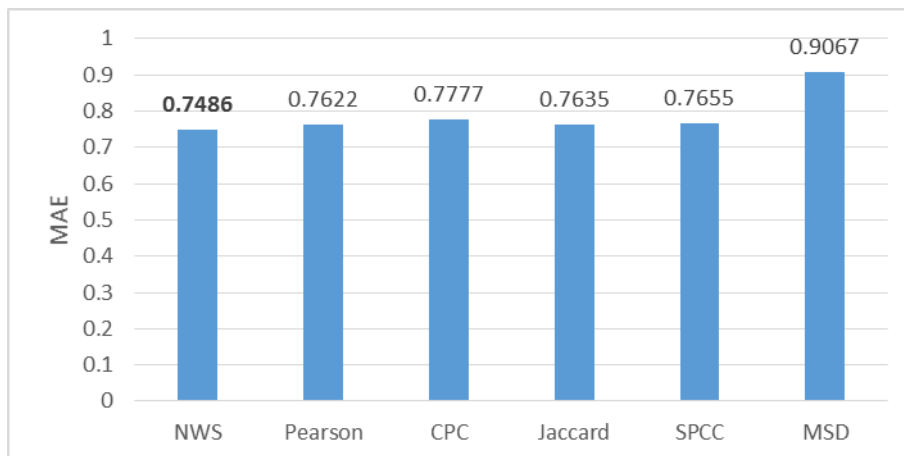
در شکل (۲) نرخ خطای سیستم توصیه‌گر پیشنهادی برای کاربران بین ۲۰ تا ۳۹ سال نسبت به سایر معیارهای شباهت به ترتیب ۱،۸ درصد، ۳،۷ درصد، ۱،۹۵ درصد، ۲،۲۱ درصد و ۱۷،۴ درصد بهبود یافته و در شکل (۳) برای کاربران بین ۴۰ تا ۶۰ سال به ترتیب ۱۳،۱ درصد، ۸ درصد، ۰،۱۱ درصد، ۱۲،۲ درصد و ۲۰،۲ درصد بهبود یافته است. دیگر نتایج آزمایشات در شکل‌های (۴) تا (۱۱)، در بیشتر معیارهای ارزیابی و شرایط، بیانگر بهبود عملکرد و کارایی سیستم توصیه‌گر پیشنهادی است.

20. Mean Absolute Error

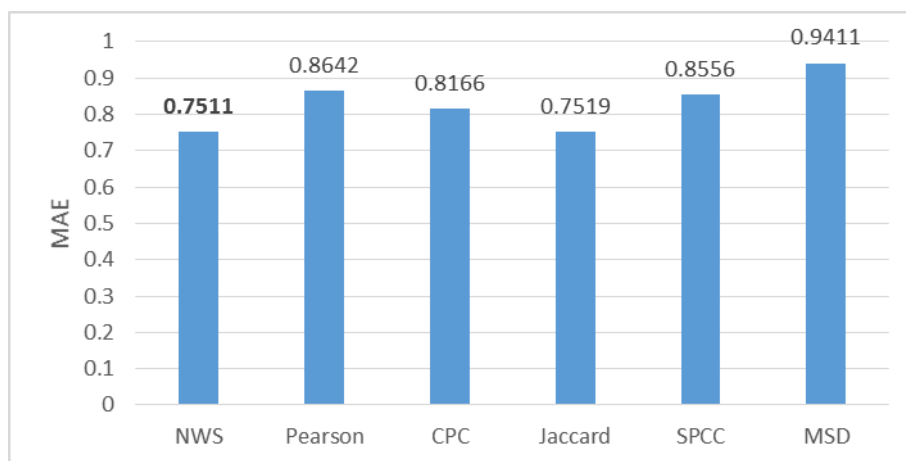
21. Accuracy

22. Precision

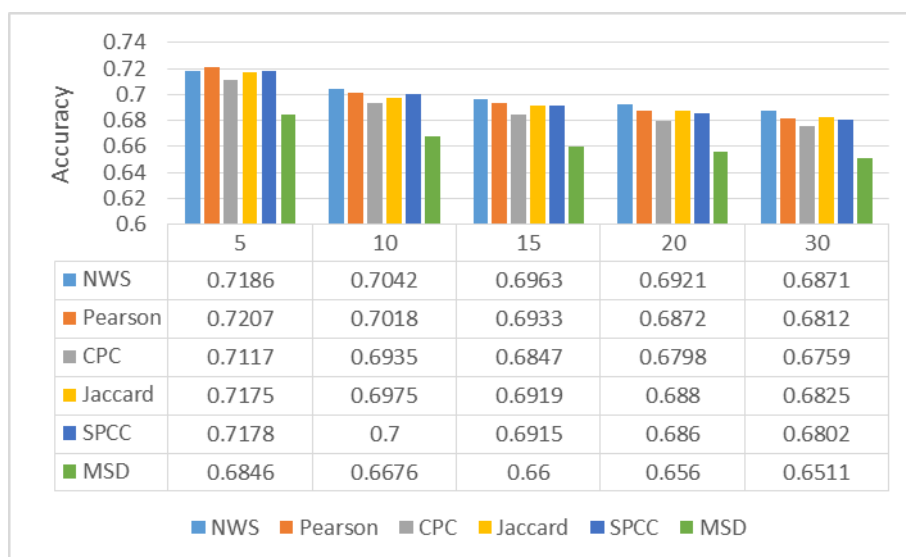
23. Recall



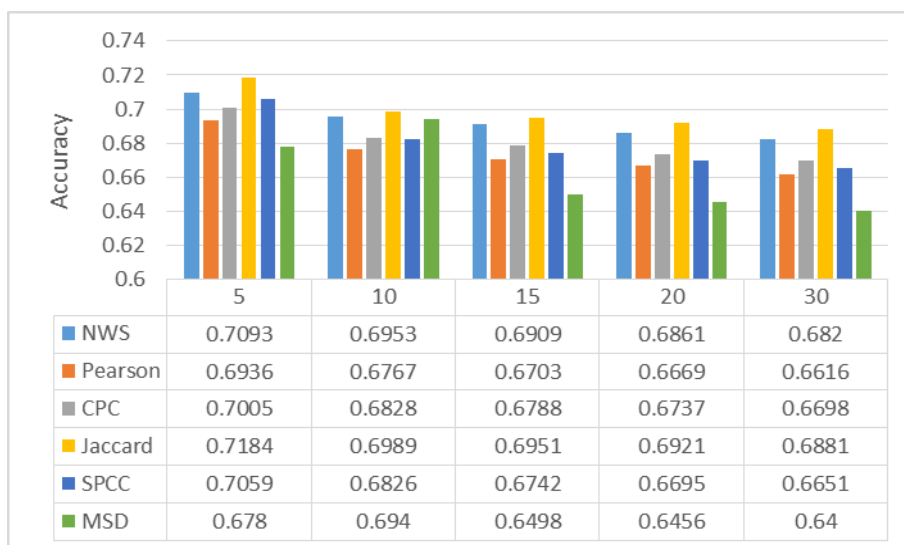
شکل ۲: میانگین مطلق خطای پیش‌بینی سیستم برای بازه‌سنی ۲۰ تا ۳۹ سال



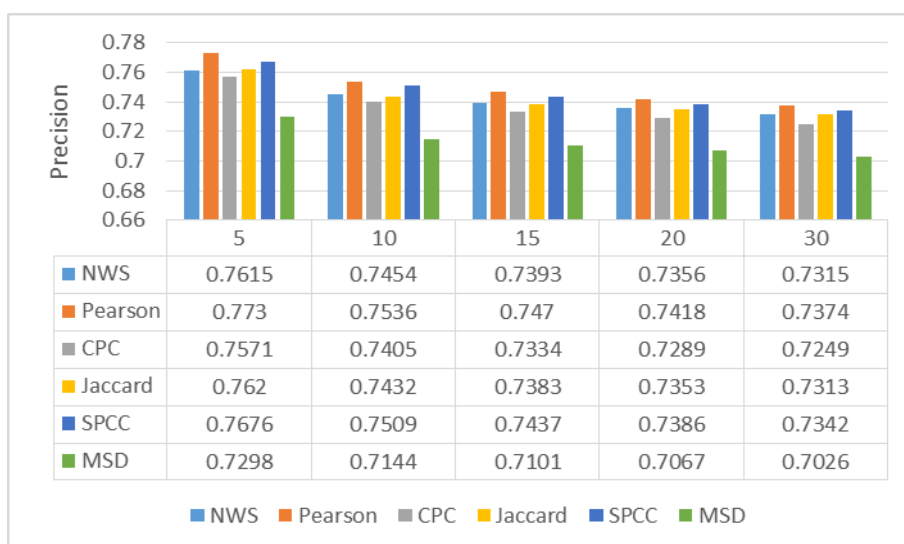
شکل ۳: میانگین مطلق خطای پیش‌بینی سیستم برای بازه‌سنی ۴۰ تا ۶۰ سال



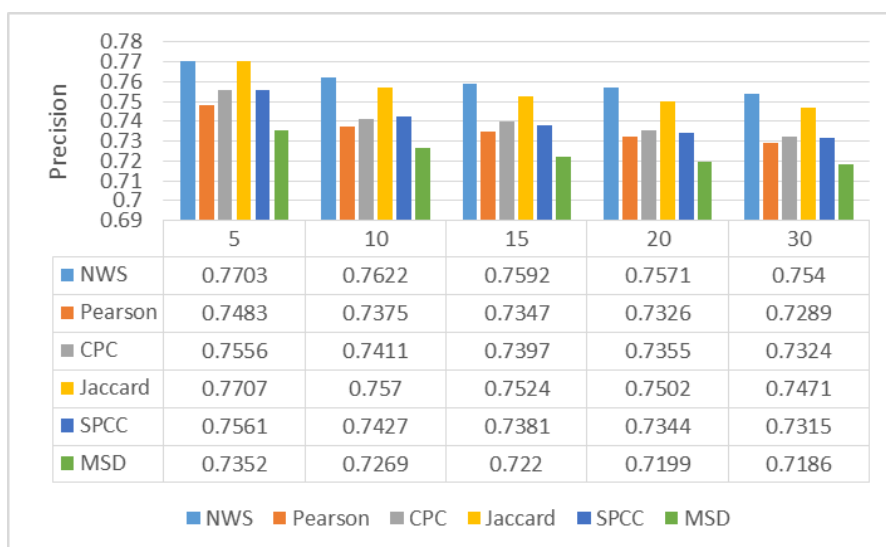
شکل ۴: دقت سیستم برای بازه‌سنی ۲۰ تا ۳۹ سال



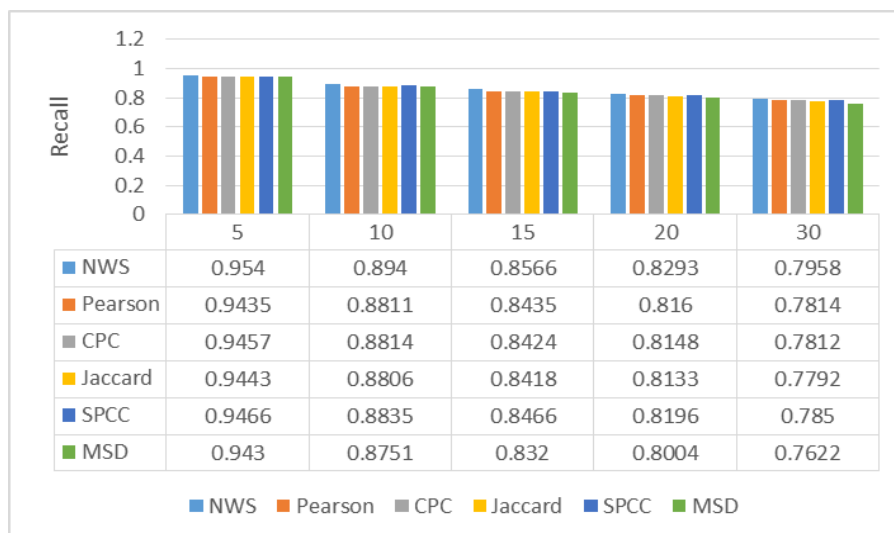
شکل ۵: دقت سیستم برای بازه‌سنی ۴۰ تا ۶۰ سال



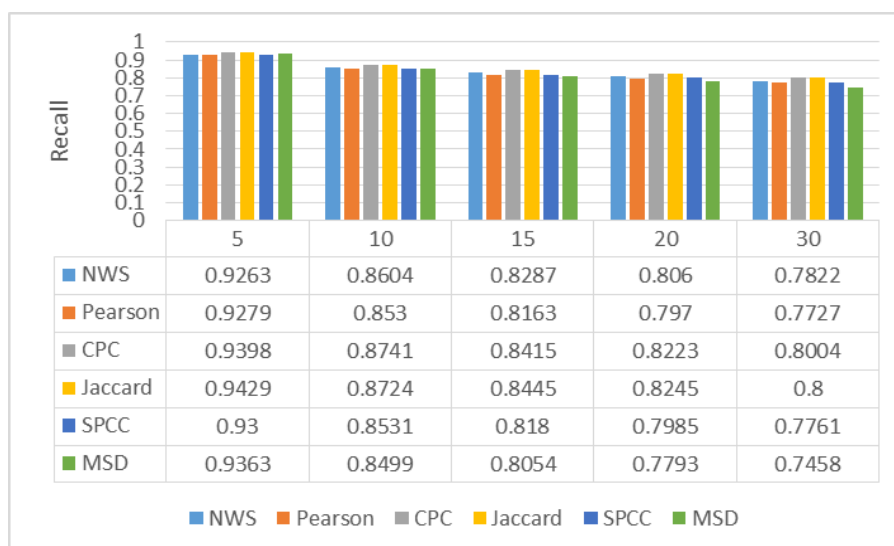
شکل ۶: صحت سیستم برای بازه‌سنی ۲۰ تا ۳۹ سال



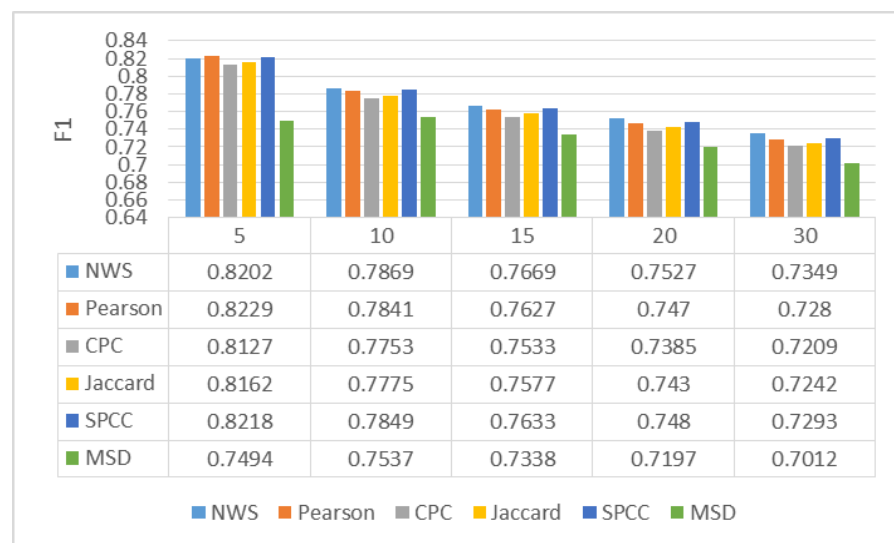
شکل ۷: صحت سیستم برای بازه‌سنی ۴۰ تا ۶۰ سال



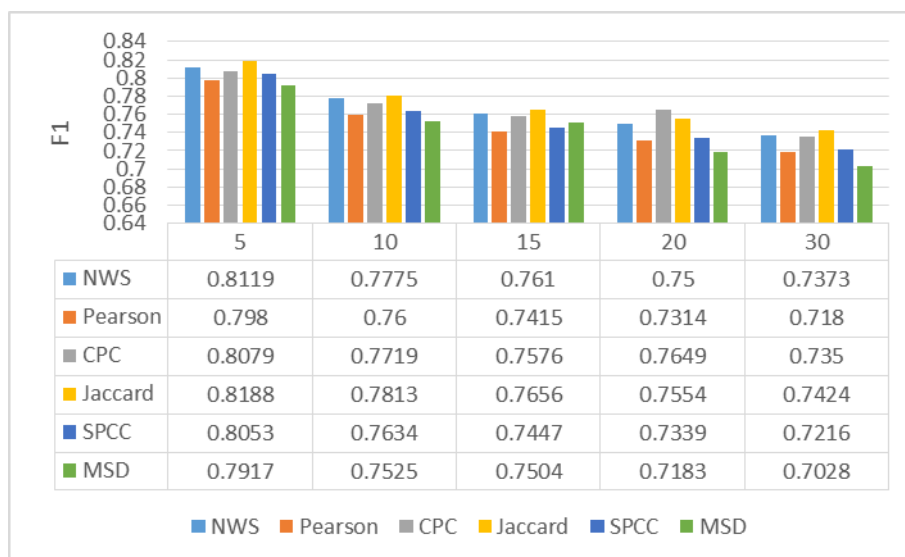
شکل ۸: بازیابی سیستم برای بازه‌سنی ۲۰ تا ۳۹ سال



شکل ۹: بازیابی سیستم برای بازه‌سنی ۴۰ تا ۶۰ سال



شکل ۱۰: معیار F1 سیستم برای بازه‌سنی ۲۰ تا ۳۹ سال



شکل ۱۱: معیار F1 سیستم برای بازه‌سنی ۴۰ تا ۶۰ سال

نتیجه‌گیری

سیستم‌های توصیه‌گر به عنوان کاربردی هوشمند و تحت وب جهت تجارت الکترونیک که امروزه توسط وبسایت‌ها ارائه می‌گردد، بسیار کارا و سودمند هستند و با پیشنهاداتش به کاربران کمک می‌کند تا تصمیم بهتری در زمان کوتاه‌تر بگیرند و سود و فروش سایت را افزایش دهد. پالایش همکارانه یکی از مهم‌ترین رویکردهای استفاده شده در سیستم‌های توصیه‌گر است. این رویکرد از مشکلات مقیاس‌پذیری، تنگی داده و خطای پیش‌بینی بالا رنج می‌برد. هدف این مقاله، طراحی یک سیستم توصیه‌گر برای شخصی‌سازی پیشنهادات براساس اطلاعات سنی کاربران به همراه یک معیارشبهت وزن‌دار جدید با نام NWS است. دسته‌بندی کاربران براساس بازه‌های سنی، باعث بهبودپذیری مساله مقیاس‌پذیری گشته و تضعیف اثر مشکل تنگی داده با بدست آوردن درجه اطمینان تشابه کاربران توسط معیارشبهت جدید NWS، باعث بهبود خطای پیش‌بینی و عملکرد سیستم می‌شود. این سیستم با تفکیک کردن کاربران به دو بازه سنی ۲۰ تا ۳۹ سال و ۴۰ تا ۶۰ سال، و اهمیت دادن به کاربرانی که تعداد فیلم‌های مشترک زیادی با کاربرفعال دارند، توانست پیشنهادات را برای کاربران فعال شخصی‌سازی کند. نتایج آزمایشات نسبت به سایر روش‌های توصیه‌گری معمول، بیانگر بهبود نرخ خطا، عملکرد و کارایی سیستم است. حداکثر نرخ خطای بهبود یافته برای کاربران بین ۲۰ تا ۳۹ سال ۱۷,۴ درصد و برای کاربران بین ۴۰ تا ۶۰ سال ۲۰,۲ درصد، می‌باشد.

سپاسگزاری

بدینوسیله نویسنده از داوران ارجمند برای نظرات ارزنده و پیشنهاداتی در راستای بهبود ارائه این مقاله، تشکر می‌کند.

منابع و مراجع

- [1] Aggarwal C. Recommender systems. Springer International Publishing, 2016.
- [۲] درزی، م.، مرادی‌منش، ز.، اصغری، ح.، "مقدمه‌ای بر سیستم‌های توصیه‌گر: الگوریتم‌ها و کاربردها". سازمان انتشارات جهاد دانشگاهی، تهران، ۱۳۸۹.
- [۳] [۳] خلجی، م.، دادخواه، چ.، FNHSM_HRS: سیستم توصیه‌گر ترکیبی مبتنی بر خوشه‌بندی فازی و معیار شباهت اکتشافی، "هفتمین کنگره مشترک سیستم‌های فازی و هوشمند ایران، بجنورد، ص ۵۶۲-۵۶۸، بهمن ۱۳۹۷.
- [4] Resnick, P., Iacovou, N., Suchak, M., Bergstrom, P. and Riedl, J. (1994, October). GroupLens: an open architecture for collaborative filtering of netnews. In Proceedings of the 1994 ACM conference on Computer supported cooperative work (pp. 175-186). ACM.
- [5] Koutrika, G., Bercovitz, B. and Garcia-Molina, H. (2009, June). FlexRecs: expressing and combining flexible recommendations. In Proceedings of the 2009 ACM SIGMOD International Conference on Management of data (pp. 745-758). ACM.
- [6] Jamali, M. and Ester, M. (2009, June). Trustwalker: a random walk model for combining trust-based and item-based recommendation. In Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 397-406). ACM.
- [7] Shardanand, U. and Maes, P. (1995, May). Social information filtering: algorithms for automating "word of mouth". In Chi (Vol. 95, pp. 210-217).
- [8] Cacheda, F., Carneiro, V., Fernández, D. and Formoso, V. (2011). Comparison of collaborative filtering algorithms: Limitations of current techniques and proposals for scalable, high-performance recommender systems. ACM Transactions on the Web (TWEB), 5(1), 2.
- [9] Goldberg D., Nichols D., Oki B. and Terry D. (1992). Using collaborative filtering to weave an information tapestry, Communications of the ACM, 35.12, 61-70.
- [۱۰] خداوردی، ن.، دادخواه، چ.، خلجی، م.، "سامانه توصیه‌گر ترکیبی مبتنی بر خوشه‌بندی و محبوبیت،" کنفرانس بین‌المللی فناوری و نوآوری در علوم، مهندسی و تکنولوژی، ایران، تهران، دانشگاه شهید بهشتی، اسفند ۱۳۹۷.
- [11] Ahn, H. J. (2008). A new similarity measure for collaborative filtering to alleviate the new user cold-starting problem. Information Sciences, 178(1), 37-51.
- [12] Liu, H., Hu, Z., Mian, A., Tian, H. and Zhu, X. (2014). A new user similarity model to improve the accuracy of collaborative filtering. Knowledge-Based Systems, 56, 156-166.
- [13] Herlocker, J., Konstan, J. A. and Riedl, J. (2002). An empirical analysis of design choices in neighborhood-based collaborative filtering algorithms. Information retrieval, 5(4), 287-310.
- [14] McLaughlin, M. R. and Herlocker, J. L. (2004, July). A collaborative filtering algorithm and evaluation metric that accurately model the user experience. In Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval (pp. 329-336). ACM.
- [15] Bellogín, A., Castells, P. and Cantador, I. (2013, May). Improving memory-based collaborative filtering by neighbour selection based on user preference overlap. In Proceedings of the 10th conference on open research areas in information retrieval (pp. 145-148). LE CENTRE DE HAUTES ETUDES INTERNATIONALES D'INFORMATIQUE DOCUMENTAIRE.
- [16] Choi, K. and Suh, Y. (2013). A new similarity function for selecting neighbors for each target item in collaborative filtering. Knowledge-Based Systems, 37, 146-153.

- [17] Javari, A., Gharibshah, J. and Jalili, M. (2014). Recommender systems based on collaborative filtering and resource allocation. *Social Network Analysis and Mining*, 4(1), 234.
- [18] Zhang, J., Lin, Y., Lin, M. and Liu, J. (2016). An effective collaborative filtering algorithm based on user preference clustering. *Applied Intelligence*, 45(2), 230-240.
- [19] Park, Y., Park, S., Jung, W. and Lee, S. G. (2015). Reversed CF: A fast collaborative filtering algorithm using a k-nearest neighbor graph. *Expert Systems with Applications*, 42(8), 4022-4028.
- [۲۰] خلجی، م.، بهبود عملکرد سیستم پیشنهاددهنده ترکیبی با استفاده از شبکه عصبی و تخصیص منابع، پایان‌نامه کارشناسی ارشد هوش مصنوعی و رباتیک، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی خواجه نصیرالدین طوسی، ایران، تهران، شهریور ۱۳۹۶.
- [21] Koochi, H., and Kiani, K. (2016). User based Collaborative Filtering using fuzzy C-means. *Measurement*, 91, 134-139.
- [22] Belacel, N., Durand, G., Leger, S. and Bouchard, C. (2018, January). Scalable Collaborative Filtering Based on Splitting-Merging Clustering Algorithm. In *International Conference on Agents and Artificial Intelligence* (pp. 290-311). Springer, Cham.
- [23] Kant, S., Mahara, T., Jain, V. K., Jain, D. K. and Sangaiah, A. K. (2018). LeaderRank based k-means clustering initialization method for collaborative filtering. *Computers & Electrical Engineering*, 69, 598-609.
- [24] Khalaji, M. and Mohammadnejad N. (2019). FCNHSMRA_HRS: Improve the performance of the movie hybrid recommender system using resource allocation approach, In *Proceedings of the 4th International Conference on Researchers in Science & Engineering & International Congress on Civil, Architecture and Urbanism in Asia*, Kasem Bundit University, Bangkok, Thailand.
- [25] Wang, Y., Feng, C., Guo, C., Chu, Y. and Hwang, J. N. (2019, January). Solving the Sparsity Problem in Recommendations via Cross-Domain Item Embedding Based on Co-Clustering. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining* (pp. 717-725). ACM.