

ارائه روشی جدید برای کشف نزدیکترین همسایگی در سیستم‌های توصیه‌گر

مبتنی بر فیلترینگ مشارکتی

مقاله پژوهشی

DOR: ۲۰,۱۰۰۱,۱,۲۷۸۳۲۵۷۰,۱۴۰۰,۲,۱,۳,۴

مهدي بازرگانی^{۱*}، زینب همایونپور^۲

۱- استادیار، گروه مهندسی کامپیوتر - دانشگاه آزاد اسلامی واحد زنجان - زنجان - ایران

۲- کارشناسی ارشد، گروه مهندسی کامپیوتر - دانشگاه آزاد اسلامی واحد زنجان - زنجان - ایران

چکیده: سیستم‌های توصیه‌گر با تحلیل و بررسی داده‌های متعلق به کاربران، یکسری آیتم‌های خاص را بر مبنای علایق به کاربران پیشنهاد می‌کنند. هدف از آنالیز داده‌های مربوط به کاربران، استخراج الگوهای هر کاربر به منظور پیش‌بینی آیتم‌ها می‌باشد. یکی از مهمترین روش‌ها در سیستم‌های توصیه‌گر، روش فیلترینگ مشارکتی است. در سیستم‌های توصیه‌گر مبتنی بر فیلترینگ مشارکتی از معیارهای شباهت جهت کشف کردن کاربران مشابه با کاربر جدید برای ارائه پیشنهاد استفاده می‌شود. از چالش‌های سیستم‌های توصیه‌گر مبتنی بر فیلترینگ مشارکتی می‌توان به فاکتورهای شباهت و تشخیص همسایگی اشاره کرد. در این مقاله از روش نزدیک‌ترین همسایه به منظور تشخیص همسایگان مشابه به کاربر جدید بر مبنای فاصله استفاده می‌کنیم. مدل پیشنهادی که برگرفته از روش کاربر-آیتم است، امتیاز قلم داده بر مبنای فاصله محاسبه می‌شود و نزدیکترین فاصله به منظور تشابه انتخاب می‌شود. در مدل پیشنهادی، تشخیص کاربران مشابه بر اساس ماتریس کاربر-آیتم توسط فاصله اقلیدسی انجام می‌شود. آزمایشات مدل پیشنهادی بر روی مجموعه داده MovieLens که شامل ۱۶۸۲ آیتم است انجام شده است. برای ارزیابی از معیارهای دقت، فراخوانی، $F1$ ، میانگین خطای مطلق و میانگین خطای مربعات ریشه استفاده شده است. میانگین خطای مطلق در مدل پیشنهادی در مقایسه با شباهت پیرسون و کسینوسی کمتر است و مقدار آن برابر با ۰/۷۳۱۵ می‌باشد. در نتیجه دقت مدل پیشنهادی در تشخیص تشابه و پیش‌بینی بیشتر است.

واژه‌های کلیدی: سیستم‌های توصیه‌گر، فیلترینگ مشارکتی، نزدیکترین همسایه، میانگین خطای مطلق

Providing A New Method to Discover the Closest Neighborhood in Recommendation Systems Based on Collaborative Filtering

Mahdi Bazergani^{1*}, Zeinab Homayounpour²

۱- Assistant Professor, Faculty of Electrical and Computer Engineering, Islamic Azad University of Zanjan, Zanjan

۲- MSc, Faculty of Electrical and Computer Engineering, Islamic Azad University of Zanjan, Zanjan

Abstract: Recommender systems, by analyzing data belonging to users, suggest a number of specific items based on interests to users. The purpose of analyzing user data is to extract each user's patterns in order to predict items. One of the most important methods in recommender systems is participatory filtering method. In collaborative filtering recommendation systems, similarity criteria are used to identify users similar to the new user to submit a proposal. One of the challenges of collaborative filtering-based recommendation systems is similarity and neighborhood detection factors. In this paper, we use the nearest neighbor method to identify similar neighbors to the new user based on distance. The proposed model, which is derived from the user-item method, calculates the data font score based on the distance, and the closest distance is selected for similarity. In the proposed model, the identification of similar users based on the user-item matrix is done by the Euclidean distance. The experiments of the proposed model were performed on the MovieLens dataset which contains ۱۶۸۲ items. Accuracy, recall, $F1$, mean absolute error and mean root square error were used for evaluation. The average absolute error in the proposed model is less than Pearson and cosine similarity and its value is equal to ۰,۷۳۱۵. As a result, the proposed model is more accurate in detecting similarity and prediction.

Keywords: Recommender Systems, Collaborative Filtering, nearest neighbor, Mean Absolute Error

تاریخ ارسال مقاله: ۱۴۰۰/۰۲/۲۸

تاریخ پذیرش مقاله: ۱۴۰۰/۰۴/۲۷

* نویسنده مسئول

۱. مقدمه

سیستم‌های توصیه‌گر یک مدل خاص از سیستم‌های پشتیبان تصمیم (سیستم‌های اطلاعاتی) می‌باشند که نقش مهمی در راهنمایی و هدایت کاربران در میان حجم عظیمی از انتخاب‌های ممکن دارند تا محصولات مناسب را برطبق اولویت‌ها و علایق خود انتخاب کنند [۱]. سیستم‌های توصیه‌گر در حوزه‌های مختلفی جهت شخصی‌سازی پیشنهادها بکار گرفته شده‌اند که از جمله آن‌ها می‌توان به کتاب، موزیک، فیلم، اخبار، مقالات و غیره اشاره کرد. سیستم‌های توصیه‌گر می‌توانند به عنوان برنامه‌هایی تعریف شوند که مناسب‌ترین قلم داده (محصول یا سرویس) را به کاربران توصیه کنند [۲]. این توصیه‌ها با پیش‌بینی سلیقه کاربر به یک نمونه خاص براساس اطلاعات مرتبط در مورد نمونه‌ها، کاربران و ارتباطات بین کاربران و نمونه‌ها ارائه می‌شود. سیستم‌های توصیه‌گر با توجه به سلیقه شخصی افراد با در نظر گرفتن انتخاب‌های قبلی وی خدماتی را به آن‌ها ارائه می‌نمایند. هدف از توسعه سیستم‌های توصیه‌گر، کاهش سربار اطلاعات به وسیله بازبانی نزدیکترین اطلاعات و خدمات از حجم انبوهی از داده‌ها و در نتیجه ارائه خدمات شخصی‌شده می‌باشد. کاربرد همه منظوره از وب باعث شده تا کاربران با حجم عظیمی از داده‌ها و اطلاعات سر و کار داشته باشند. کاربران در مواجهه با این حجم اطلاعات دچار سردرگمی در انتخاب اطلاعات مورد نظر خود می‌شوند. سیستم‌های توصیه‌گر ابزار بسیار مهم و کارآمدی در وب سایت‌های تجارت الکترونیک محسوب می‌شوند و به کاربران در کشف قلم داده‌مورد علاقه‌شان کمک می‌کنند. هدف اصلی سیستم‌های توصیه‌گر تولید پیشنهادهای دقیق می‌باشد. این سیستم‌ها می‌توانند پیشنهادهایی را براساس پروفایل کاربران و یا اولویت‌های قبلی آن‌ها ارائه کنند و یا اینکه بر انتخاب‌های سایر افراد به عنوان داور تکیه کنند. مهم‌ترین ویژگی یک سیستم توصیه‌گر قابلیت پیش‌بینی اولویت‌ها و سلیقه‌های یک کاربر به وسیله تحلیل رفتار آن و رفتار سایر کاربرها جهت تولید پیشنهادهای شخصی شده است [۳].

فیلترینگ مشارکتی معمول‌ترین روش مورد استفاده در سیستم‌های توصیه‌گر است. این روش مبتنی بر ارزیابی کاربرانی است که علاقه‌های مشابهی دارند و ایده اصلی این سیستم‌ها این است که کاربرانی که آیت‌های مشابهی را در گذشته انتخاب کرده‌اند احتمالاً ارجحیت‌های مشابهی دارند. سیستم‌های فیلترینگ مشارکتی از داده‌های قدیمی مرتبط با اولویت‌های کاربر یا رفتار آن استفاده می‌کنند تا بتوانند رفتار کاربران جدید را پیش‌بینی کنند. هدف از فیلترینگ مشارکتی پیش‌بینی رتبه مربوط به اقلام‌های رتبه‌بندی نشده در ماتریس کاربر-آیتم می‌باشد [۴]. در فیلترینگ مشارکتی فرض بر این است که اگر دو کاربر رتبه مشابهی را به یک آیتم (کتاب، فیلم، موزیک و...) بدهند یا رفتارهای مشابهی (خرید کردن، مطالعه کردن، تماشا کردن و غیره) را داشته باشند، آنها می‌توانند اقلام‌های مشابه را رتبه‌بندی کنند. مبتنی بر این فرض، فیلترینگ مشارکتی بر مبنای نرخ تشابه کاربران و اقلام‌ها به پیش‌بینی رتبه اقلام‌های رتبه‌بندی نشده می‌پردازد. در این سیستم‌ها

یک مدل یا سند بر مبنای ارزیابی مجموعه‌ای از کاربران ایجاد می‌شود در واقع کاربران را گروه‌بندی می‌کند، هر گروه شامل کاربرها و اقلام‌هایی است که شبیه به هم می‌باشند. فیلترینگ مشارکتی دارای ضعف‌هایی زیر است: در این روش سیستم باید پیشنهاداتی را بر مبنای رتبه‌ی کاربر به سایر اقلام‌ها به وی ارائه دهد. از معایب دیگر این روش می‌توان به اضافه شدن قلم داده جدید به مجموعه‌ی اقلام‌ها اشاره کرد. تا زمانیکه یک قلم داده دارای رتبه‌ای نباشد، آن قلم داده به کاربر جدید پیشنهاد نمی‌گردد. به علاوه مشکل پراکندگی در این الگوریتم وجود دارد. یعنی تعداد رتبه‌هایی که از گذشته بدست آمده بسیار کمتر از تعداد مورد نیاز برای پیشگویی است [۵].

فیلترینگ مشارکتی به طور وسیعی نسبت به روش‌های دیگر در سیستم‌های توصیه‌گر مورد استفاده قرار گرفته است. این سیستم‌ها پیش‌بینی را از طریق سابقه همسایگان (کاربران یا آیتم‌ها) محاسبه می‌کنند و انتخاب صحیح نزدیک‌ترین همسایه که بیشترین شباهت را به کاربر هدف داشته باشند تأثیر زیادی بر دقت پیش‌بینی‌ها خواهد داشت. در روش پیشنهادی، الگوریتم نزدیکترین همسایه بر مبنای تکنیک‌های تشخیص فاصله در جهت کسب نتایج بهتر استفاده می‌شود.

در سیستم‌های توصیه‌گر بیشتر از روش فیلترینگ مشارکتی استفاده شده است. مشکلات مهم این سیستم‌ها، پراکندگی، مقیاس‌پذیری و هم‌معنایی می‌باشد که در دقت و کیفیت پیشنهاددهی تأثیر منفی می‌گذارد. با توجه به این که افزایش دقت پیش‌بینی در سیستم‌های توصیه‌گر فیلترینگ مشارکتی نقش اساسی در بسیاری از زمینه‌های زندگی از جمله تجارت، پزشکی و... ایفا می‌کند، تاکنون روش‌ها و الگوریتم‌های متعددی جهت رفع این مشکلات ارائه شده است، اما هنوز چالش مهم این سیستم‌ها، افزایش دقت آن می‌باشد. بنابراین، نیاز به ارائه‌ی روش‌های هوشمند جهت بهبود دقت این سیستم‌ها لازم و اساسی می‌باشد و در این تحقیق در راستای رفع این نیاز از روش مبتنی بر مدل استفاده می‌شود. الگوریتم‌های مبتنی بر مدل، پیش‌بینی رتبه‌بندی را با استفاده از روش‌های یادگیری ماشین و آماری، برای یادگیری یک مدل از داده‌ها انجام می‌دهند. جنبه جدید بودن این تحقیق در استفاده از الگوریتم نزدیکترین همسایه و بکارگیری روش‌های مختلف فاصله برای تشخیص همسایگان می‌باشد.

۲. مطالعات پیشین

در طرح [۶] پروین و همکارانش از الگوریتم بهینه‌سازی کلونی مورچه به منظور پیش‌بینی رأی‌های از دست رفته مبتنی بر فیلترینگ مشارکتی استفاده کرده‌اند. روش پیشنهادی شامل سه مرحله اصلی است: در مرحله اول، کاربران با در نظر گرفتن مقدار ارزش و روابط اعتماد رتبه‌بندی می‌شوند. سپس، در مرحله دوم، الگوریتم بهینه‌سازی کلونی مورچه برای تعیین مقادیر مناسب وزن به کاربران مورد استفاده قرار می‌گیرد. یک مجموعه از کاربران مشابه در مرحله سوم فیلتر می‌شوند تا در پیش‌بینی رتبه‌های ناشناخته برای کاربر هدف مورد استفاده قرار گیرند. به عبارت

دیگر، برای سرعت بخشیدن به شناسایی کاربران مشابه، روش پیشنهادی ابتدا بخش اکثریت کاربران غیرمشابه را فیلتر می‌کند و سپس الگوریتم بهینه‌سازی کلونی مورچه را فقط در یک مجموعه کاهش یافته از کاربران اجرا می‌کند تا آنها را وزن‌دهی کند. در هر تکرار مورچه‌ها ارزیابی می‌شوند و بر مبنای انتخاب‌شان، بهترین مقدار برای موارد ناشناخته لحاظ می‌گردد. شبیه‌سازی بر روی سه مجموعه داده E-Film-Trust، Ciao و pinions انجام شده است. نتایج نشان داده است که الگوریتم بهینه‌سازی کلونی مورچه از مقدار خطای کمتری در مقایسه با تکنیک‌های دیگر بهره‌مند است.

در طرح [۷] کانت و همکارانش از الگوریتم خوشه‌بندی k-means به منظور گروه‌بندی آیتم‌های مشابه استفاده کرده‌اند که هدف آنها یافتن مراکز خوشه بهینه به منظور گروه‌بندی آیتم‌ها در گروه‌های مشابه است. الگوریتم خوشه‌بندی k-means بر روی ماتریس داده‌های کاربر مورد استفاده قرار می‌گیرد تا محدوده مورد نظر را تعیین کند. الگوریتم k-means خوشه‌بندی قلم داده را بر اساس اندازه‌گیری فاصله تعریف شده بین دو آیتم مختلف پیدا می‌کند. الگوریتم خوشه‌بندی k-means گروه‌هایی از قلم داده را شناسایی می‌کند به طوری که فقط متعلق به یک خوشه باشند. شبیه‌سازی بر روی مجموعه داده MovieLens ۱۰۰k انجام شده است. نتایج نشان می‌دهد که الگوریتم k-means بهبود داده شده در مقایسه با معمولی خطای کمتری دارد.

در طرح [۸] واسید و همکارش یک مدل بر مبنای الگوریتم بهینه‌سازی اجتماع ذرات و فازی برای فیلترینگ مشارکتی پیشنهاد داده‌اند. در طرح آنها از الگوریتم بهینه‌سازی اجتماع ذرات برای وزندهی به ویژگی‌ها و از مجموعه فازی برای انتخاب ویژگی‌های موثر استفاده شده است. در این کار، با استفاده از الگوریتم بهینه‌سازی ذرات، اولویت‌ها بر مبنای رتبه ایجاد می‌شوند. پس از پیدا کردن وزن مناسب برای ویژگی رتبه، بهترین تشابه به کاربر فعال نسب داده می‌شود. شبیه‌سازی بر روی مجموعه داده MovieLens انجام شده است. نتایج نشان داده که روش پیشنهادی در مقایسه با روش‌های دیگر از درصد خطای کمتری بهره‌مند است.

در طرح [۹] جو و همکارش از الگوریتم زنبور مصنوعی برای خوشه‌بندی K-means استفاده کرده‌اند. در فرایند خوشه‌بندی از الگوریتم زنبور مصنوعی برای رفع مشکل ناشی از خوشه‌بندی K-means استفاده شده است. سپس شباهت بین کاربران یک خوشه، توسط شباهت کسینوسی بهبود داده شده، محاسبه گردیده است. شبیه‌سازی مدل ترکیبی بر روی مجموعه داده MovieLens انجام شده است.

در طرح [۱۰] آر و بوستانسی به منظور کاهش خطای پیش‌بینی از الگوریتم ژنتیک و معیار شباهت بردار کسینوسی استفاده کرده‌اند. از الگوریتم ژنتیک برای کاربران همسایه بسیار مشابه وزندهی شده و با استفاده از معیار شباهت و وزن‌های به دست آمده، دقت در انتخاب همسایگان افزایش یافته است و با استفاده از جهش ژنتیکی به کاربران جدید امکان اضافه شدن به لیست کاربران مشابه را داده و به این ترتیب

باعث بهبود مشکل شروع سرد شده است. شبیه‌سازی مدل ترکیبی بر روی مجموعه داده MovieLens انجام شده است. نتایج نشان داده که مدل پیشنهادی در مقایسه با مدل‌های سنتی از درصد خطای کمتر و افزایش دقت برخوردار بوده است.

در طرح [۱۱] کوهی و کیانی از روش خوشه‌بندی C میانگین فازی یا FCM برای فیلترینگ مشارکتی مبتنی بر کاربر استفاده کرده‌اند و عملکرد آن در برابر روش‌های خوشه‌بندی مختلف مورد بررسی قرار گرفته است. رویه این مدل به پنج بخش تقسیم می‌شود: (۱) آماده‌سازی داده‌ها؛ خواندن داده و پاکسازی داده‌ها. (۲) خوشه‌بندی کاربر؛ الگوریتم خوشه‌بندی FCM برای تولید P بخش از کاربران، اعمال شده است. این خوشه‌بندی به عنوان پایه‌ای برای پیدا کردن خوشه‌های دیگر، استفاده شده است. (۳) خوشه‌بندی پویای فازی: خوشه‌های فازی در قالب‌های زمانی مختلف، برای یک کاربر فعال، پیدا شده و درجه پویایی عضویت، محاسبه شده است. (۴) انتخاب همسایه: همسایگی برای یک کاربر داده شده بر اساس خوشه‌بندی فازی، تعیین می‌شود. (۵) توصیه: رتبه‌بندی کاربر فعال از آیتم رتبه‌بندی نشده بر اساس رتبه‌بندی‌های همسایه‌ها، پیش‌بینی می‌شود. داده‌های ورودی در این روش، از یک ماتریس کاربر به رتبه‌بندی، وارد شده‌اند. مجموعه داده مووی لنز (MovieLens) برای مقایسه الگوریتم‌های خوشه‌بندی مختلف مورد استفاده قرار گرفته است. و روش‌ها از نظر دقت و صحت و بازخوانی ارزیابی شده‌اند. نتایج شبیه‌سازی نشان داده که ترکیبی از روش خوشه‌بندی فازی با دیفازی کردن مرکز ثقل و ضریب همبستگی پیرسون منجر به توصیه‌های بهتری به نسبت سایر روش‌ها شده است.

در طرح [۱۲] مدل خوشه‌بندی بهبودیافته-میانگین به منظور کاهش مشکل پیچیدگی محاسباتی پیشنهاد شده است که کاربران را به خوشه‌هایی همگن و منسجم تقسیم‌بندی کرده است. خوشه‌بندی بر مبنای معیار شباهت فاصله انجام شده است که کاربران درون هر خوشه دارای شباهت بالا و شباهت بین خوشه‌ها در سطح پایینی می‌باشد. مدل پیشنهادی به صورت بازگشتی داده‌های حجیم نظرات را به دو زیرخوشه تقسیم کرده است. درخت حاصل از این بخش‌بندی، یک درخت دودویی غیرمتعادل بوده است که هر برگ آن دارای یک ماتریس شباهت و در هر گره داخلی، مراکز نظرات زیردرخت‌ها ثبت شده است. به منظور پیش‌بینی نظر، ابتدا خوشه‌ای که کاربر به آن تعلق داشته است مشخص شده و سپس از نظرات موجود در همان خوشه برای تشخیص نظر نزدیک به کاربر فعال استفاده شده است. کارایی مدل پیشنهادی به تعداد خوشه‌ها و اندازه خوشه‌ها وابسته بوده است. در زمان اجرا، هر چقدر اندازه خوشه‌ها کوچک‌تر باشد، بار محاسباتی کمتر ولی از طرفی دیگر کیفیت پیش‌بینی نظرات کاهش پیدا کرده است.

در [۱۳] از الگوریتم خفاش برای محاسبه وزن آیتم‌ها (ویژگی‌ها) استفاده شده است تا همسایگی بهتری برای کاربر فعال پیدا شود. این تکنیک برای وزندهی به آیتم‌ها با استفاده از الگوریتم خفاش در دستیابی به

توصیه‌های بهتر کمک کرده است. عملکرد این سیستم با عملکرد سیستم مبتنی بر کلونی زنبورهای مصنوعی نیز مقایسه شده است. نتایج نشان داد که الگوریتم خفاش از نظر میانگین خطای مطلق و معیار $F1$ در مقایسه با الگوریتم کلونی زنبور مصنوعی در حدود ۶,۹ درصد بهتر بوده است.

در [۱۴]، یک مدل فیلترینگ مشارکتی مبتنی بر یادگیری عمیق ارائه شده است. این مدل شامل دو بخش است: بخش اول این مدل مربوط به ویژگی‌های کاربران و آیت‌ها است و از آنها به عنوان ورودی در رتبه‌بندی شبکه عصبی عمیق استفاده می‌شود، که رتبه‌بندی‌ها را پیش‌بینی می‌کند. این رتبه‌بندی‌ها ورودی بخش دوم را تشکیل می‌دهند که یک رتبه‌بندی مبتنی بر شبکه عصبی عمیق است و برای پیش‌بینی رتبه‌بندی کلی استفاده می‌شود. آزمایشات بر روی یک مجموعه داده واقعی نشان داده است که مدل پیشنهادی از روش‌های دیگر بهتر عمل کرده است و مقدار خطای کمتری داشته است.

در [۱۵] یک سیستم توصیف فیلم مبتنی بر مدل ترکیبی که از k -means و ژنتیک تشکیل شده است برای پارتیشن‌بندی فضای کاربران ارائه شده است. این روش از تجزیه و تحلیل مؤلفه‌های اصلی داده برای کاهش تراکم فضای جمعیت فیلم استفاده می‌کند که می‌تواند پیچیدگی محاسبات را کاهش دهد. الگوریتم K -Means ابتدا K عضو (که K تعداد خوشه‌ها است) را به صورت تصادفی از میان N عضو انتخاب می‌نماید و آنها را به عنوان مراکز خوشه‌ها در نظر می‌گیرد. سپس $N-K$ عضو باقیمانده به نزدیکترین خوشه تخصیص می‌یابند. بعد از تخصیص همه اعضا مجدداً مراکز خوشه‌ها محاسبه شده و اعضا با توجه به میزان نزدیکی (شباهت) به یکی از خوشه‌ها تخصیص می‌یابند و این کار تا زمانی که مراکز خوشه‌ها ثابت بمانند، ادامه می‌یابد. با تکرار همین روال می‌توان در هر تکرار با میانگین‌گیری از داده‌ها مراکز جدیدی برای آنها محاسبه کرد و مجدداً داده‌ها را به خوشه‌های جدید نسبت داد. ارزش رتبه‌بندی تخمین زده شده برای یک فیلم بدون رتبه‌بندی توسط کاربر Ua طبق معادله (۱) تعریف شده است.

(۱)

$$P_{Ua,item} = \bar{R}_u + \frac{\sum_{y \in C_x} sim(Ua, y) \times (R_{y,i} - \bar{R}_y)}{\sum_{y \in C_x} (|sim(Ua, y)|)}$$

در معادله (۱) میانگین رتبه‌بندی بر مبنای Ua ، مجموعه‌ای از همسایه‌ها متعلق به خوشه Ua ، پارامتر \bar{R}_y میانگین رتبه‌بندی بر مبنای امتیازهای دریافت شده از همسایگان Ua ، پارامتر $sim(Ua, y)$ تابع تشابه بر مبنای معیار پیرسون برای درجه تشابه دو کاربر. نتایج آزمایش بر روی مجموعه داده Movielens نشان داده است که روش پیشنهادی توانسته از لحاظ دقت، عملکرد بالایی را ارائه دهد و در مقایسه با روش‌های موجود، پیشنهادی قابل اطمینان‌تر و شخصی‌تری تولید کند.

در [۱۶] از یک روش فیلترینگ گروهی برای ارائه توصیه به کاربران استفاده شده است. در اکثر روش‌های فیلترینگ گروهی از متدهایی مانند

ضریب پیرسون و روش‌های مبتنی بر کسینوس برای بدست آوردن میزان شباهت دو کاربر استفاده می‌شود. اما در این تحقیق فاصله کاربران با استفاده از هفت روش مختلف (فاصله اقلیدسی، فاصله چبیشف، فاصله گاور، فاصله سورنسن، فاصله کانبرا، فاصله لورنتزین و فاصله سی‌تی بلاک) اندازه‌گیری و محاسبه شده است. با توجه به فاصله محاسبه شده از دو الگوریتم k نزدیک‌ترین همسایه و الگوریتم k -means برای مشخص شدن کاربران مشابه به کاربر هدف استفاده شده است. نتایج نشان داده است که دقت این مدل تقریباً ۶۰ درصد است و میانگین خطا چیزی در حدودی ۲۵ درصد بوده است.

در [۱۷] یک سیستم پیشنهادی برای فیلترینگ مشارکتی بر مبنای ترکیب میانگین‌گیری و الگوریتم جلبک مصنوعی پیشنهاد شده است. برای یافتن شباهت بین دو کاربر از ضریب همبستگی پیرسون چندسطحی پیشرفته استفاده شده است. علاوه بر این، رتبه به فیلم‌هایی که هنوز کاربر به احتمال زیاد رتبه‌ای برای آنها اختصاص نداده است پیدا می‌شود. سیستم پیشنهادی در مقایسه با روش‌های دیگر، موفق به ارائه توصیه‌هایی با کیفیت و دقت بهتر بوده است. سیستم پیشنهادی در چهار مجموعه داده واقعی، آزمایش و ارزیابی شده است: Movielens، ۱۰۰، ۰۰۰، Epinion و Jester. سیستم پیشنهادی برای هر چهار مجموعه داده توصیه‌های بهتری ارائه داده است. با استفاده از معیارهای ارزیابی از جمله میانگین خطای مطلق (MAE)، دقت و فراخوان، کارایی سیستم برآورد شده است.

در [۱۸] یک سیستم توصیه‌گیر مبتنی بر الگوریتم ژنتیک^۱ ارائه شده است که به اطلاعات معنایی و داده‌های رتبه‌بندی قبلی بستگی دارد. ایده اصلی این تحقیق در ارزیابی لیست توصیه‌های احتمالی به جای ارزیابی آیت‌ها و سپس تشکیل لیست توصیه‌ها است. BLIGA از الگوریتم ژنتیک برای یافتن بهترین لیست موارد به کاربر فعال استفاده می‌کند. بنابراین، هر کروموزم لیست پیشنهادی کاندیداها را نشان می‌دهد. BLIGA سلسله مراتبی، کروموزم‌ها را با استفاده از سه تابع برازندگی ارزیابی کرده است. تابع اول از اطلاعات معنایی آیت‌ها استفاده کرده تا قدرت شباهت معنایی بین آیت‌ها را تخمین بزند. تابع دوم شباهت سطح رضایت کاربران را تخمین زده است. تابع سوم به رتبه‌بندی پیش‌بینی شده بستگی داشته است تا بهترین لیست توصیه‌ها را انتخاب نماید. نتایج نشان داده است که برتری BLIGA و توانایی آن برای دستیابی به پیش‌بینی دقیق‌تر از روش‌های دیگر بوده است.

در [۱۹] یک چارچوب جدید بر مبنای فازی به منظور مسئله شروع سرد پیشنهاد شده است. در این روش برای پیشنهاد آیت‌ها به کاربران جدید از قوانین فازی استفاده شده است. در معادله (۲) p و q مجموعه اعداد فازی هستند که بر مبنای آیت‌های قبلی تنظیم شده‌اند.

$$S_{u_j, u_k}(i_p) = 1 - \frac{1}{\sqrt{5}} \left(\sum_{h=1}^5 (q_{jp}^h - q_{kp}^h)^2 \right)^{1/2} \quad (2)$$

مدل فازی شامل مراحل استنتاج فازی، مراحل تصمیم‌گیری و توصیه می‌باشد. نتایج نشان داده است که مدل پیشنهادی میانگین خطای پایین و درصد دقت بالایی بالایی داشته است.

در [۲۰] به منظور افزایش دقت پیش‌بینی بهبود سیستم‌های فیلترینگ مشارکتی، از طرحی مبتنی بر الگوریتم ژنتیک استفاده شده است. در این طرح به طور انتخابی مقادیر اندازه‌گیری شباهت با بالاترین مقدار انتخاب می‌شوند و پروفایل کاربری در داخل یک کروموزوم، کد می‌شود. مدل به سه فاز تقسیم می‌شود: فاز اول: فاز جمع‌آوری اطلاعات کاربری است که در آن برای یک فرد جدید، سیستم از او می‌خواهد که به منظور جمع‌آوری ویژگی‌ها شخصی و همچنین رتبه‌بندی‌ها از یک سری از فیلم‌ها، ثبت نام کند. فاز دوم ایجاد ویژگی‌های پروفایل کاربر است که برای آن دسته از فیلم‌هایی که کاربر رتبه بندی کرده است. فاز سوم که فاز توصیه‌گر الگوریتم ژنتیک نامیده شده است، دارای توابع مرتبط با الگوریتم ژنتیک است که توصیه مناسب را از طریق انتخاب و وزندهی ویژگی‌ها، دنبال می‌کند و سپس به انجام توصیه می‌پردازد. در این روش برای اندازه‌گیری شباهت نیز از ضریب همبستگی پیرسون استفاده شده است. برای ارزیابی این روش، آزمایش‌هایی بر روی مجموعه داده MovieLens بر مبنای معیار دقت انجام شده است و مشخص شده است که این روش با مقدار دقت بین ۸۱/۲۸ تا ۸۱/۷۷ درصد نسبت به روش پیرسون پایه با مقدار دقت ۶۹/۲ درصد، دقت بالاتری را نشان داده است. مزیت این روش این است که همیشه یک جواب نسبتاً خوب تولید می‌کند و دقت نسبت به روش‌های پیشین افزایش یافته است. عیب این روش این است که نیاز به حافظه و زمان زیادی دارد و لزوماً بهترین جواب را به دست نمی‌آورد.

۳. مدل پیشنهادی

فیلترینگ مشارکتی، یکی از محبوب‌ترین و کارآمدترین مدل‌ها برای تعریف پیشنهاد بالقوه به کاربران جدید است. دو نوع الگوریتم فیلترینگ مشارکتی وجود دارد: سیستم توصیه‌گر مبتنی بر کاربر، که آیتم‌ها را با رتبه‌بندی شباهت آیتم‌های مشترک بین کاربران قبلی به کاربران جدید پیشنهاد می‌کند. فیلترینگ مشارکتی مبتنی بر آیتم، با در نظر گرفتن رتبه آیتم‌ها و تاریخچه آنها، آیتم را به کاربران جدید توصیه می‌کند. اگرچه فیلترینگ مشارکتی با موفقیت در بسیاری از زمینه‌های تجاری مورد استفاده قرار گرفته‌اند، اما چندین اشکال اصلی فیلترینگ مشارکتی، به ویژه کمبود داده‌های رتبه‌بندی، چالش جدی برای صحت و جهانی بودن این روش‌ها ایجاد کرده است. به طور خاص، بیشترین تعداد رتبه‌بندی برای هر کاربر در بسیاری از آیتم‌ها وجود ندارد و عملکرد فیلترینگ مشارکتی به همراه افزایش تعداد آیتم بدون رتبه در مجموعه داده‌های آموزشی، کاهش می‌یابد.

الگوریتم نزدیک‌ترین همسایه در بیشتر موارد به منظور دسته‌بندی آیتم‌های مشابه به کار می‌رود، هرچند که می‌توان از آن برای تخمین و پیش‌بینی نیز استفاده نمود. در الگوریتم نزدیک‌ترین همسایه، آیتم‌ها بر مبنای فاصله کشف می‌شوند. معیار فاصله برای تعیین تشابه بین آیتم‌ها استفاده می‌شود. اگر دو بردار q و p داشته باشیم از فاصله اقلیدسی طبق معادله (۳) برای بدست آوردن فاصله بین دو آیتم q_i و p_i استفاده می‌شود.

$$d(q, p) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2} \quad (3)$$

در الگوریتم نزدیک‌ترین همسایه، باید در بین آیتم‌ها، آیتمی انتخاب شود که برای آیتم بدون امتیاز کمترین فاصله را داشته باشد. در نتیجه، پس از آنکه فاصله اقلیدسی بین آیتم‌ها محاسبه شد، با مرتب‌سازی عناصر بر حسب فاصله اقلیدسی، از میان k همسایه، نمونه‌ای که از میان k همسایه دارای کمترین مقدار است به نمونه ناشناخته داده می‌شود.

در سیستم‌های توصیه‌گر تاریخچه علاقه‌مندی‌های کاربر به آیتم‌های خریداری شده در ماتریس رتبه ذخیره می‌شود. سطرهای این ماتریس نمایانگر کاربران و ستون‌های این ماتریس نشانگر آیتم‌های خریداری شده توسط کاربران هستند. هر عنصر این ماتریس، نشان‌دهنده‌ی میزان علاقه یک کاربر خاص به آیتم خاص می‌باشد. میزان علاقه یک کاربر به یک آیتم خاص به صورت صریح و یا ضمنی می‌تواند بدست بیاید، یعنی کاربر به آیتم‌هایی که انتخاب کرده یا به صورت صریح رتبه داده است یا به صورت خودکار به کالای خریداری شده توسط سیستم رتبه اختصاص داده شده است. در این ماتریس هر کاربر تنها به تعداد محدودی آیتم انتخاب شده رتبه داده است و میزان علاقه‌مندی خود را به آن آیتم‌های مورد نظر مشخص نموده است و نظر آن نسبت به بسیاری از آیتم‌ها مشخص نیست. در شکل (۱) مراحل مدل پیشنهادی نشان داده شده است.

الگوریتم‌های فیلترینگ مشارکتی انتخاب مناسب یک تابع شباهت می‌باشد. زیرا دقت و عملکرد توصیه را تحت تاثیر قرار می‌دهد. الگوریتم نزدیک‌ترین همسایه در فیلترینگ مشارکتی مبتنی بر حافظه براساس معیار فاصله، شباهت را محاسبه می‌کند. این معیار بر مبنای مقادیری که توسط کاربران به هر آیت‌م داده شده است تعریف می‌شود.

۲.۳. گام پیش‌بینی

بعد از تولید ماتریس مجاورتی با استفاده از امتیازهای متعلق به هر آیت‌م که توسط کاربران همسایه امتیاز داده شده‌اند، امتیاز هر آیت‌م برای کاربر جدید محاسبه می‌شود و در قالب پیش‌بینی ارائه می‌گردد. هدف از این مرحله پیش‌بینی امتیازهای هر آیت‌م بدون امتیاز برای کاربران جدید با استفاده از فیلترینگ مشارکتی مبتنی بر آیت‌م است. امتیازهای جدید پیش‌بینی شده به عنوان امتیازی برای آیت‌م‌های بدون رتبه در نظر گرفته می‌شوند و به عنوان نتایج نهایی ثبت می‌شوند. تعریف پیش‌بینی طبق معادله (۵) انجام می‌شود. پس از به دست آوردن یک مجموعه از همسایه‌ها و تعریف فیلم‌های امتیاز داده شده توسط همسایه‌ها می‌توان لیستی از توصیه‌ها را به کاربر جدید پیشنهاد داد که توسط وی امتیازبندی نشده‌اند. همچنین رأی آن فیلم‌ها باید پیش‌بینی گردد که کاربر جدید با انتخاب آیت‌م‌های دلخواه به ثبت رأی و پیش‌بینی کمک می‌کند. رأی برای فیلم u برای کاربر جدید u به منظور پیش‌بینی طبق معادله (۵) تعریف می‌شود.

$$\text{prediction}_{u,i} = \bar{v}_u + k \sum_{j=1}^n d(u,j)(v_{j,i} - \bar{v}_j) \quad (5)$$

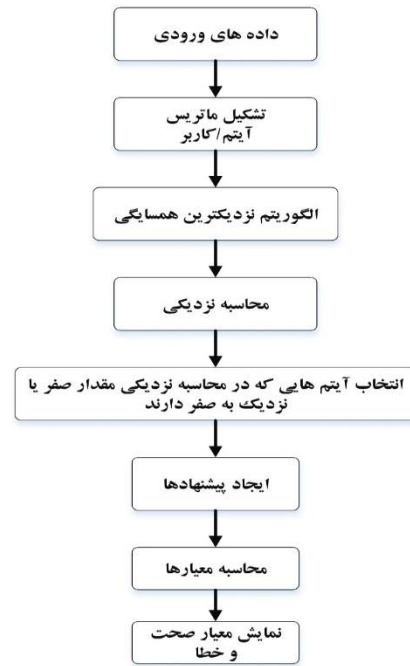
پارامتر \bar{v}_u میانگین امتیاز از کاربر u ، پارامتر k فاکتور نرمالسازی برای محدوده بهینه که بین صفر و یک است. پارامتر d فاصله اقلیدسی، پارامتر n تعداد همسایه‌ها، پارامتر $v_{j,i}$ مقدار امتیاز واقعی از همسایه j ام برای آیت‌م u ام است.

۳.۳. معیارهای ارزیابی

باید به منظور کارایی و درستی مدل، تحلیل و ارزیابی سیستم‌های توصیه‌گر بر مبنای معیارهای مختلف انجام شود. برای ارزیابی سیستم‌های توصیه‌گر معیارهای گوناگونی در منابع مختلف تعریف شده است. از جمله معمولی‌ترین معیارهای ارزیابی سیستم‌های توصیه‌گر معیارهای دقت، فراخوانی و معیار $F1$ می‌باشد. دقت، یکی از مهمترین اهداف طراحی در سیستم‌های توصیه‌گر می‌باشد. در ادامه معیارهای مورد نظر را توضیح خواهیم داد و در این تحقیق از معیارهای بیان شده برای مقایسه عملکرد و کارایی مدل پیشنهادی استفاده خواهیم کرد.

۱.۳.۲. میانگین خطای مطلق

به منظور اندازه‌گیری دقت نتایج سیستم‌های توصیه‌گر پیشنهادی، معمولاً از رایج‌ترین معیارهای پیش‌بینی خطا استفاده می‌شود که از



شکل ۱: مراحل مدل پیشنهادی بر مبنای ساختار فلوجارتی

۱.۳.۳. مرحله شباهت

در این مرحله با ارزیابی روابط ماتریس کاربر-آیت‌م باید زیر مجموعه‌ای از نزدیک‌ترین کاربران به کاربر جدید (شبهه‌ترین کاربران به کاربر جدید) تشخیص داده می‌شوند. در این مرحله، از فاصله اقلیدسی بر مبنای امتیازات کاربر قبلی و کاربر جدید برای بدست آوردن میزان شباهت استفاده می‌شود. به منظور ارزیابی میزان شباهت مقدار امتیاز آیت‌م‌ها محاسبه می‌شوند اگر مقدار بدست آمده نزدیک به صفر باشد آنگاه شباهت بیشتر است. برای مثال اگر کاربر جدید به یکی از آیت‌م‌ها امتیاز ۴ داده باشد، این مقدار با مقدار امتیازهای قبلی محاسبه می‌شود و مقدار فاصله اقلیدسی آیت‌م‌هایی که با آیت‌م جدید برابر با صفر باشد به عنوان تشابه انتخاب می‌شوند.

اگر در محاسبه فاصله بین کاربران میزان تشابه وجود نداشته باشد آنگاه همسایه‌های دیگر بررسی می‌شوند. داده‌های ورودی در فیلترینگ مشارکتی معمولاً به عنوان یک ماتریس ارزیابی $R(m, n)$ از کاربران $m \times n$ بیان می‌شود، به طوری که سطر m نشانگر کاربران M است و ستون n به معنای آیت‌م N است. R_{ij} مقدار رتبه آیت‌م j توسط کاربر i است. تعریف رابطه بین کاربران-آیت‌م‌ها به صورت معادله (۴) است. هر کاربر به آیت‌م‌ها یک رتبه داده است که به صورت عدد ذخیره شده است.

$$R_{mn} = \begin{bmatrix} R_{11} & \dots & R_{1n} \\ \vdots & \ddots & \vdots \\ R_{m1} & \dots & R_{mn} \end{bmatrix} \quad (4)$$

در یک سیستم توصیه‌گر با n کاربر و m آیت‌م، مجموعه کاربران را به صورت U و مجموعه آیت‌م‌ها را با I نشان می‌دهیم. ماتریس کاربر-آیت‌م از ورودی‌های اصلی یک سیستم توصیه‌گر می‌باشد. این ماتریس شامل رتبه‌ی کاربران به آیت‌م‌های موجود در سیستم است. یک فاکتور مهم در

جمله آن‌ها میانگین خطای مطلق یا همان MAE^x و معیارهای وابسته به آن مثل میانگین خطای مربعات ریشه یا همان $RMSE^x$ می‌باشد. این دو معیار، خطای پیش‌بینی کاربر را محاسبه می‌کنند. میانگین خطای مطلق، میانگین انحراف بین امتیاز پیش‌بینی شده و امتیاز واقعی کاربر را محاسبه می‌کند. از جمله مزیت‌های MAE این است که نحوه‌ی محاسبه‌ی آن ساده و قابل فهم بوده و همچنین دارای ویژگی‌های آماری می‌باشد. هر چه مقدار MAE کمتر باشد، پیش‌بینی دقیق‌تر خواهد بود. فاکتور MAE طبق معادله (۶) تعریف می‌شود.

$$MAE = \frac{1}{N} \sum_{u,i} |p_{u,i} - r_{u,i}| \quad (6)$$

که در آن $r_{u,i}$ به ترتیب امتیاز پیش‌بینی شده و امتیاز واقعی کاربر u برای آیت i و N تعداد کل امتیازات به مجموعه‌ی آیت‌ها می‌باشد. همچنین ریشه میانگین خطای مطلق به صورت معادله (۷) تعریف می‌شود.

$$RMSE = \sqrt{\frac{1}{N} \sum_{u,i} (p_{u,i} - r_{u,i})^2} \quad (7)$$

۲.۳.۳. معیارهای صحت

هدف از معیارهای صحت و جامعیت محاسبه‌ی میزان درست یا نادرست بودن توصیه‌های صورت گرفته می‌باشد. معیار صحت درستی پیشنهادها را اندازه‌گیری می‌کند. در واقع این معیار بیان می‌کند که از مجموعه پیشنهادهایی که به کاربر داده می‌شود چند درصد درست است بنابراین مقدار بالای این معیار نشان‌دهنده‌ی تعداد اشتباهات کمتر می‌باشد معیار مورد نظر به صورت معادله (۸) تعریف می‌شود.

$$Precision = \frac{TP}{TP+FP} \quad (8)$$

در معادله (۸) پارامتر TP تعداد آیت‌های پیشنهاد شده مورد علاقه کاربران جدید می‌باشد در حالیکه FP تعداد آیت‌های پیشنهادی است که کاربران جدید به آن علاقه ندارند. مجموع TP و FP تعداد کل پیشنهادات داده شده به کاربران جدید می‌باشد. معیار فراخوانی نشان می‌دهد که چند درصد از آیت‌های مورد علاقه کاربر به وی پیشنهاد شده است. این معیار به صورت معادله (۹) محاسبه می‌شود. که در آن FN تعداد کاربران جدیدی است که به آیت‌ها علاقه دارند ولی به آن‌ها پیشنهاد نشده است. و مجموع TN و FN مجموعه‌ای از آیت‌های مورد علاقه کاربران جدید است که باید به آن‌ها پیشنهاد شود.

$$Recall = \frac{TP}{TP + FN} \quad (9)$$

در معادله (۹) بررسی می‌شود که آیا کاربر از پیشنهاداتی که سیستم به وی داده است، راضی می‌باشد یا خیر. به عبارتی دیگر درصد رضایت کاربر از آیت‌های پیشنهاد شده به وی توسط معیار فراخوانی، ارزیابی می‌شود. دو معیار دقت و فراخوانی در تقابل با یکدیگر تغییر می‌کنند. ویژگی دو

معیار دقت و فراخوانی با هم متفاوت است. غالباً افزایش تعداد کل آیت‌ها باعث افزایش در نتایج فراخوانی و کاهش مقدار دقت می‌شود. برای یافتن توازن مناسب بین مقادیر دقت و فراخوانی معیار دیگری به نام معیار $F1$ که معمولاً با نام معیار F اندازه‌گیری بیان می‌شود تعریف شده است. این معیار مقادیر دقت و فراخوانی را در یک معیار ترکیب می‌کند و به هر یک از آن‌ها وزن‌های مساوی تخصیص می‌دهد. در واقع این معیار، معیارهای دقت و فراخوانی را به وسیله‌ی میانگین هارمونیک آن‌ها با هم ترکیب می‌کند. مقدار $F1$ به شدت به سمت مقادیر کمتر از دقت و فراخوانی تمایل دارد. هر چه مقدار $F1$ به ۱۰۰ درصد نزدیک‌تر باشد حکایت از دقت بیشتر خروجی و نزدیک بودن آن به خواست کاربر دارد.

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall} \quad (10)$$

۴. ارزیابی و نتایج

برای آزمایش مدل پیشنهادی از نرم‌افزار متلب ۲۰۱۶ استفاده شده است. نرم‌افزار متلب که توسط کمپانی مس‌ورک توسعه یافته، به عنوان یکی از پیشرفته‌ترین نرم‌افزارهای علمی شناخته می‌شود که می‌تواند برای محاسبه عددی داده‌ها، توسعه الگوریتم‌ها و تجزیه و تحلیل داده‌ها مورد استفاده قرار گیرد. حجم مجموعه داده‌ها بسیار بالا است و نیاز به محاسبات بالایی دارد. برای راه‌اندازی برنامه نیاز به مواردی همانند خواندن فایل ورودی، تشکیل ماتریس کاربر-آیت، محاسبه متوسط امتیازهایی که یک کاربر به آیت‌ها داده است، و به دست آوردن لیست آیت‌هایی که کاربر به آنها امتیاز داده است. در این پژوهش، آزمایش‌ها بر روی سیستم عامل ویندوز ۸/۱ و با پردازنده‌ی ۶۴ بیتی و دو هسته‌ای، هر هسته ۲/۵ گیگاهرتز و حافظه ۴ گیگابایتی اجرا شده است.

۴.۱.۴. ارزیابی و نتایج: معیارهای صحت

در این مقاله جهت ارزیابی نتایج شبیه‌سازی از مجموعه داده MovieLens استفاده شده است [۲۱]. بر اساس بررسی‌های انجام شده، مجموعه داده MovieLens شامل ۱۰۰۰۰۰ رکورد رتبه‌ای برای ۱۶۸۲ آیت است که توسط ۹۴۳ کاربر ثبت شده است و مقیاس رتبه‌دهی به صورت اعداد ۱ تا ۵ است. نرخ پراکندگی برابر با ۹۵/۶۹ درصد و حداکثر تعداد رای‌های کاربران برابر با ۷۳۷ و میانگین تعداد رای‌های کاربران برابر با ۱۰۶ است. در ابتدا، تاثیر تعداد همسایگی بر معیارهای صحت بررسی شده است. ابتدا مقادیر دقت، فراخوانی و $F1$ بر روی مجموعه داده MovieLens با استفاده از تعداد همسایه‌ها (۵، ۱۰، ۱۵، ۲۰، ۲۵، ۳۰، ۴۰ و ۵۰) محاسبه شده است. نتایج در جدول (۱) نشان داده شده است که بهترین دقت متعلق به مدل پیشنهادی است.

دقت مدل پیشنهادی با ۱۰ همسایه برابر با ۰/۵۳۱۴ و با ۵۰ همسایه برابر با ۰/۴۰۱۱ است. فراخوانی مدل پیشنهادی با ۱۰ همسایه برابر با ۰/۵۴۱۷ و با ۵۰ همسایه برابر با ۰/۴۱۲۶ است. نتایج نشان می‌دهد اگر تعداد همسایه‌ها بیشتر شود درصد دقت و فراخوانی کمتر می‌شود، به

دلیل اینکه فاصله بین آیتم‌ها بیشتر می‌شود و تشخیص بهترین آیتم برای همسایگی ممکن است با اشتباه صورت بگیرد.

پیرسون با ۱۰ همسایه برابر با ۰/۴۳۶۴ و با ۵۰ همسایه برابر با ۰/۲۸۵۲ است. همچنین دقت شباهت کسینوسی با ۱۰ همسایه برابر با ۰/۴۴۱۳ و با ۵۰ همسایه برابر با ۰/۳۱۹۶ است. فراخوانی شباهت کسینوسی با ۱۰ همسایه برابر با ۰/۴۴۷۲ و با ۵۰ همسایه برابر با ۰/۳۲۰۸ است.

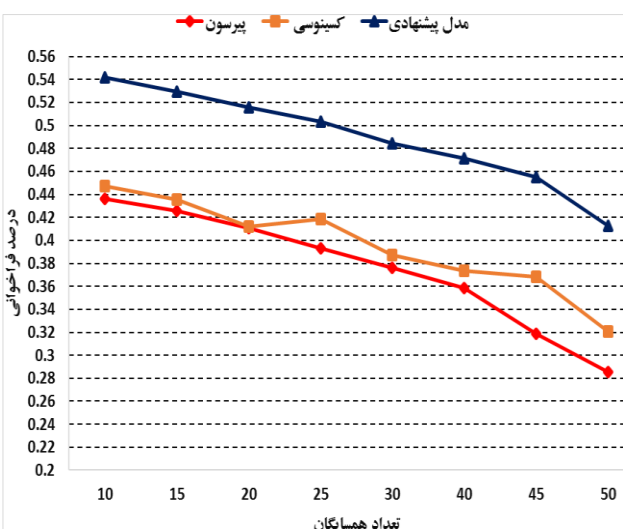
نتایج جدول (۱) نشان می‌دهد که دقت شباهت پیرسون با ۱۰ همسایه برابر با ۰/۴۳۴۱ و با ۵۰ همسایه برابر با ۰/۲۸۱۴ است. فراخوانی شباهت

جدول ۱: ارزیابی صحت در مدل پیشنهادی برحسب تعداد همسایه‌ها

مدل‌ها	معیارها	۱۰	۱۵	۲۰	۲۵	۳۰	۴۰	۴۵	۵۰
شباهت پیرسون	دقت	۰/۴۳۴۱	۰/۴۲۱۲	۰/۴۱۸۵	۰/۳۹۶۲	۰/۳۷۳۱	۰/۳۵۲۱	۰/۳۰۱۴	۰/۲۸۱۴
	فراخوانی	۰/۴۳۶۴	۰/۴۲۵۸	۰/۴۱۱۰	۰/۳۹۳۱	۰/۳۷۶۱	۰/۳۵۸۴	۰/۳۱۸۶	۰/۲۸۵۲
	F1	۰/۴۴۰۲	۰/۴۲۶۵	۰/۴۱۵۴	۰/۳۹۵۲	۰/۳۷۲۸	۰/۳۵۴۷	۰/۳۰۹۸	۰/۲۸۴۱
شباهت کسینوسی	دقت	۰/۴۴۱۳	۰/۴۳۹۷	۰/۴۲۸۹	۰/۴۱۵۲	۰/۳۹۱۶	۰/۳۸۴۶	۰/۳۵۴۷	۰/۳۱۹۶
	فراخوانی	۰/۴۴۷۲	۰/۴۳۵۶	۰/۴۱۲۳	۰/۴۱۸۶	۰/۳۸۷۴	۰/۳۷۳۲	۰/۳۶۸۴	۰/۳۲۰۸
	F1	۰/۴۴۴۲	۰/۴۳۶۵	۰/۴۱۹۵	۰/۴۱۷۲	۰/۳۸۹۵	۰/۳۷۸۵	۰/۳۵۹۶	۰/۳۱۲۵
مدل پیشنهادی	دقت	۰/۵۳۱۴	۰/۵۲۴۹	۰/۵۰۲۱	۰/۴۹۴۵	۰/۴۸۸۴	۰/۴۶۷۱	۰/۴۵۹۴	۰/۴۰۱۱
	فراخوانی	۰/۵۴۱۷	۰/۵۲۹۴	۰/۵۱۵۶	۰/۵۰۳۲	۰/۴۸۴۷	۰/۴۷۱۷	۰/۴۵۴۸	۰/۴۱۲۶
	F1	۰/۵۳۶۵	۰/۵۲۴۹	۰/۵۰۳۵	۰/۴۹۷۵	۰/۴۸۳۵	۰/۴۶۷۳	۰/۴۵۳۵	۰/۴۰۶۸

است که اگر تعداد آیتم‌های پیش‌بینی زیاد باشد آنگاه تشخیص کاربران مشابه برای کاربر جدید به دلیل میانگین رای‌های پایین ممکن نیست و مقدار فراخوانی کمتر می‌شود. اما درصد فراخوانی در شباهت پیرسون و کسینوسی در همان اوایل در مقایسه با مدل پیشنهادی کمتر است. این دو مدل برای محاسبه تشابه از محاسبه فاصله مستقیم بین کاربران قبلی و کاربر جدید استفاده نمی‌کنند و به دلیل استفاده از فاکتورهای مختلف از آیتم مشابه دور می‌شوند و تشخیص خطا بیشتر می‌شود.

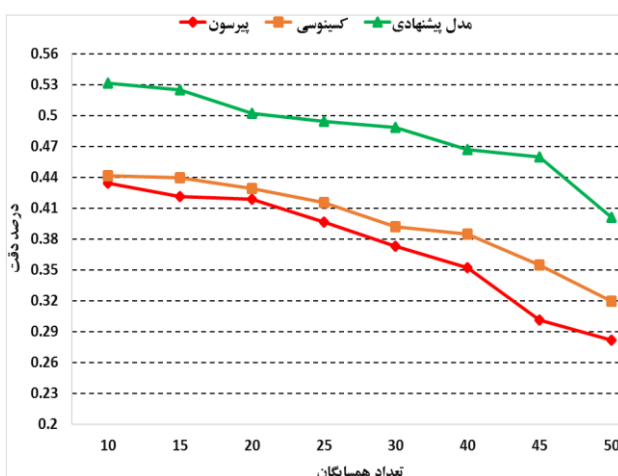
در شکل (۲) نمودار مقایسه مدل پیشنهادی با مدل‌های پیرسون و کسینوسی بر مبنای معیار دقت نشان داده شده است. نتایج در شکل (۲) نشان می‌دهد که بهترین دقت را مدل پیشنهادی و با تعداد ۱۰ همسایه ارائه داده است. معیار دقت در پیرسون و کسینوسی کمتر است و این بیانگر این است که مدل پیشنهادی مقدار خطای کمتری داشته است و همسایه را بهتر از لحاظ تشابه تشخیص داده است.



شکل ۳: نمودار مقایسه مدل پیشنهادی با مدل‌های پیرسون و

کسینوسی بر مبنای معیار فراخوانی

در شکل (۴) نمودار مقایسه مدل پیشنهادی با مدل‌های پیرسون و کسینوسی بر مبنای معیار F1 نشان داده شده است. نتایج شکل (۴) نشان

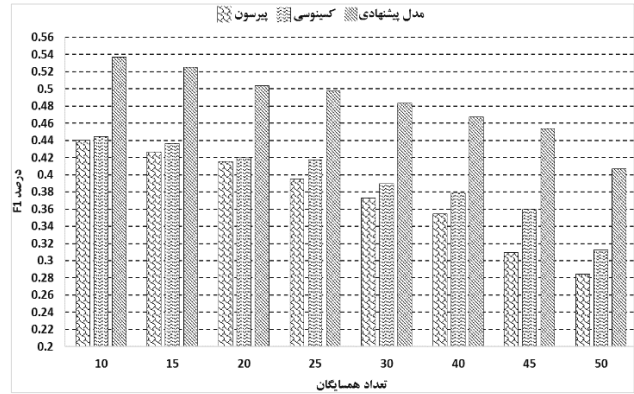


شکل ۲: نمودار مقایسه مدل پیشنهادی با مدل‌های پیرسون و

کسینوسی بر مبنای معیار دقت

در شکل (۳) نمودار مقایسه مدل پیشنهادی با مدل‌های پیرسون و کسینوسی بر مبنای معیار فراخوانی نشان داده شده است. نمودار شکل (۳) نشان می‌دهد که با افزایش تعداد همسایه‌ها در مدل پیشنهادی و مدل‌های دیگر، درصد فراخوانی کاهش پیدا می‌کند. این بدان دلیل

می‌دهد که معیار F_1 در مدل پیشنهادی در مقایسه با شباهت پیرسون و کسینوسی بیشتر است. معیار F_1 در مدل پیشنهادی در مقایسه با شباهت پیرسون و شباهت کسینوسی به ترتیب 0.1227 و 0.0943 درصد بهبود یافته است.



شکل ۴: نمودار مقایسه مدل پیشنهادی با مدل‌های پیرسون و کسینوسی بر مبنای معیار F_1

۲.۴. ارزیابی و نتایج: معیارهای خطا

جدول (۲) نشان می‌دهد که با افزایش تعداد همسایه‌ها مقدار MAE و RMSE کمتر شده است. با افزایش تعداد همسایه‌ها، مقدار MAE در مدل پیشنهادی کاهش زیادی داشته است که بیانگر افزایش قابل توجه در صحت پیش‌بینی است. نتایج نشان می‌دهد که مقدار MAE و RMSE در مدل پیشنهادی برای حالت ۵۰ همسایه برابر با 0.7315 و 0.9521 می‌باشد. مقدار MAE و RMSE در شباهت پیرسون برای حالت ۵۰ همسایه برابر با 0.9716 و 0.9241 می‌باشد. مقدار MAE و RMSE در شباهت کسینوسی برای حالت ۵۰ همسایه برابر با 0.9148 و 0.9681 می‌باشد. به طور کلی همانطور که در جدول (۲) نشان داده شده است مدل پیشنهادی توانسته است به مقادیر خطای کمتری در مقایسه با روش‌های دیگر دست یابد. مدل پیشنهادی مقدار خطای MAE را در مقایسه با پیرسون در حدود 0.1926 کمتر کرده است.

جدول ۲: ارزیابی خطا در مدل پیشنهادی بر حسب تعداد همسایه‌ها

مدل‌ها	معیارها	۱۰	۱۵	۲۰	۲۵	۳۰	۴۰	۴۵	۵۰
شباهت پیرسون	MAE	0.9574	0.9528	0.9468	0.9407	0.9348	0.9315	0.9217	0.9241
	RMSE	1.05	1.03	0.9993	0.9971	0.9874	0.9758	0.9734	0.9716
شباهت کسینوسی	MAE	0.9519	0.9476	0.9411	0.9394	0.9314	0.9288	0.9261	0.9148
	RMSE	0.9979	0.9958	0.9886	0.9862	0.9814	0.9778	0.9705	0.9681
مدل پیشنهادی	MAE	0.8841	0.8621	0.8514	0.8376	0.8149	0.7935	0.7614	0.7315
	RMSE	0.9843	0.9822	0.9785	0.9739	0.9691	0.9648	0.9579	0.9521

کاربر به آن آیت‌ها محسوب می‌شود. لذا در فرم کلی، اینگونه مسائل را می‌توان مسئله پیش‌بینی امتیاز برای آیت‌هایی در نظر گرفت که هنوز توسط کاربران دیده نشده‌اند. فیلترینگ مشارکتی، به طور معمول یک ماتریس رای که شامل نظرات کاربران نسبت به آیت‌ها است و پیش‌بینی پیشنهادها بر اساس علائق کاربران مشابه تولید و ارائه می‌شوند. در مدل پیشنهادی، برای پیدا کردن شباهت میان دو کاربر، تابعی جدید بر مبنای فاصله اقلیدسی پیشنهاد شد. تشخیص این شباهت، بر اساس نزدیکی امتیازهای دو کاربر عمل می‌کند. نزدیکی بین دو رای می‌تواند تاثیر زیادی در نتیجه نهایی داشته باشد و موجب کاهش خطا می‌شود. مدل پیشنهادی بر روی مجموعه داده MovieLens ارزیابی شد و توسط معیارهای صحت و خطا مورد سنجش و تست قرار گرفت. درصد معیارهای صحت نشان دادند که مدل پیشنهادی دارای درصد دقت بهتری در مقایسه با پیرسون و کسینوسی است و همچنین مقدار خطا در مدل پیشنهادی کمتر بود.

۵. نتیجه‌گیری

در حوزه سیستم‌های توصیه‌گر، تعدادی کاربر وجود دارد که نظرات خود را نسبت به تعدادی آیت اعلام می‌نمایند. با فرض اینکه یک سیستم توصیه‌گر دارای m کاربر و n آیت باشد و رای‌های کاربران به آیت‌ها در محدوده r_{min}, \dots, r_{max} باشد، که در آن r_{min} و r_{max} به ترتیب بیانگر کوچک‌ترین و بزرگ‌ترین مقدار رای‌ها می‌باشند. از نماد U برای نمایش مجموعه کاربران و I برای نمایش مجموعه آیت‌ها استفاده می‌شود. نظر هر کاربر نسبت به یک آیت رای نامیده می‌شود که معمولاً به صورت یک سه‌گانه $(u, i, r_{u,i})$ نشان داده می‌شود که در آن u معرف یک کاربر، i معرف یک آیت و $r_{u,i}$ نشان‌دهنده ارزش رای داده شده توسط کاربر u به آیت i می‌باشد.

مسئله مبتنی بر ساختار امتیازدهی کاربران به مجموعه‌ای از آیت‌ها را سیستم‌های توصیه‌گر می‌گویند. امتیاز داده شده هر کاربر به هر یک از آیت‌ها، نشان‌دهنده نرخ علاقه کاربر به آن آیت‌ها و تحت عنوان رای

۲۰۰۱: Data Warehousing and Knowledge Discovery, pp. ۱۴۱-۱۵۱, ۲۰۰۱.

- [۱۳] S. Yadav, Vikesh, Shreyam, S. Nagpal, An Improved Collaborative Filtering Based Recommender System using Bat Algorithm, *Procedia Computer Science*, Vol. ۱۳۲, ۱۷۹۵-۱۸۰۳, ۲۰۱۸.
- [۱۴] J.N. Ruskin, V. Fuster, Abstracts of original contributions: ۴۲nd Annual Scientific Session, *Journal of the American College of Cardiology*, Vol. ۲۱, Issue ۲, pp. ۱۷a-۴۸۸a, ۱۹۹۳
- [۱۵] Z. Wang, X. Yu, N. Feng, Z. Wang, an improved collaborative movie recommendation system using computational intelligence, *Journal of Visual Languages and Computing*, Vol. ۲۵, pp. ۶۶۷-۶۷۵, ۲۰۱۴.
- [۱۶] G.M. Dakhel, M. Mahdavi, Providing an Effective Collaborative Filtering Algorithm Based on Distance Measures and Neighbors' Voting, *International Journal of Computer Information Systems and Industrial Management Applications*, Vol. ۵, pp. ۵۲۴-۵۳۱, ۲۰۱۳.
- [۱۷] R. Katarya, O.P. Verma, Effectual recommendations using artificial algae algorithm and fuzzy c-mean, *Swarm and Evolutionary Computation*, Vol. ۳۶, pp. ۵۲-۶۱, ۲۰۱۷.
- [۱۸] B. Alhijawi, Y. Kilani, A collaborative filtering recommender system using genetic algorithm, *Information Processing and Management*, Vol. ۵۷, pp. ۱-۲۱, ۲۰۲۰.
- [۱۹] LC. Cheng, HA. Wang, A fuzzy recommender system based on the integration of subjective preferences and objective information, *Applied Soft Computing*, Vol. ۱۸, pp. ۲۹۰-۳۰۱, ۲۰۱۴.
- [۲۰] Y. Ho, S. Fong, Z. Yan, On Improving GA-based Collaborative Filtering for Online Recommender, pp. ۱-۷, ۲۰۰۷. <https://grouplens.org/datasets/movielens/۱۰۰k/>

زیرنویس

۱. Genetic-based Recommender System (BLIGA)
۲. Mean Absolute Error (MAE)
۳. Root Mean Square Error (RMSE)

برای کارهای آینده، می‌توان از خوشه‌بندی فازی به منظور تشخیص کاربران مشابه استفاده کرد. یکی از روش‌هایی که برای کنار هم قرار دادن کاربران مشابه استفاده می‌شود، خوشه‌بندی است. با استفاده از خوشه‌بندی می‌توان کاربران مشابه به هم را در یک خوشه و کاربرانی که رفتار شبیه به هم ندارند را در خوشه‌های متفاوت قرار داد. در خوشه‌بندی فازی، هر کاربر می‌تواند مربوط به چندین خوشه باشد ولی درجه یا میزان ارجحیت آن کاربر به هر خوشه، بستگی به میزان شباهت آن کاربر با مرکز خوشه مورد نظر دارد. در این حالت تعداد زیادی از کاربران مشابه با درجه‌ای مشخص، به کاربر جدید تشابه دارند.

۶. مراجع

- [۱] M. Duma, B. Twala, Optimising latent features using artificial immune system in collaborative filtering for recommender systems, *Applied Soft Computing*, Volume ۷۱, pp. ۱۸۳-۱۹۸, ۲۰۱۸.
- [۲] S. Jiang, S.C. Fang, Q. An, J.E. Lavery, A sub-one quasi-norm-based similarity measure for collaborative filtering in recommender systems, *Information Sciences*, Vol. ۴۸۷, pp. ۱۴۲-۱۵۵, ۲۰۱۹.
- [۳] D. Valcarce, A. Landin, J. Parapar, A. Barreiro, Collaborative filtering embeddings for memory-based recommender systems, *Engineering Applications of Artificial Intelligence*, Vol. ۸۵, pp. ۳۴۷-۳۵۶, ۲۰۱۹.
- [۴] F. Zhang, T. Gong, V.E. Lee, G. Zhao, Guangzhi Qu, Fast algorithms to evaluate collaborative filtering recommender systems, *Knowledge-Based Systems*, Vol. ۹۶, pp. ۹۶-۱۰۳, ۲۰۱۶.
- [۵] T. Mohammadpour, A.M. Bidgoli, R. Enayatifar, H.H.S. Javadi, Efficient clustering in collaborative filtering recommender system: Hybrid method based on genetic algorithm and gravitational emulation local search algorithm, *Genomics*, In press, corrected proof, Available online ۳ January ۲۰۱۹
- [۶] H. Parvin, P. Moradi, S. Esmaili, TCFACO: Trust-aware collaborative filtering method based on ant colony optimization, *Expert Systems with Applications*, Vol. ۱۱۸, pp. ۱۵۲-۱۶۸, ۲۰۱۹.
- [۷] S. Kant, T. Mahara, V.K. Jain, D.K. Jain, A.K. Sangaiah, LeaderRank based k-means clustering initialization method for collaborative filtering, *Computers & Electrical Engineering*, Vol. ۶۹, pp. ۵۹۸-۶۰۹, ۲۰۱۸.
- [۸] M. Wasid, V. Kant, A Particle Swarm Approach to Collaborative Filtering based Recommender Systems through Fuzzy Features, *Procedia Computer Science*, Vol. ۵۴, pp. ۴۴۰-۴۴۸, ۲۰۱۵.
- [۹] C. Ju and C. Xu, A New Collaborative Recommendation Approach Based on Users Clustering Using Artificial Bee Colony Algorithm, *Hindawi Publishing Corporation, The Scientific World Journal*, Article ID ۸۶۹۶۵۸, pp. ۱-۹, ۲۰۱۳.
- [۱۰] Y. Ar, E. Bostanci, A genetic algorithm solution to the collaborative filtering problem, *Expert Systems with Applications*, Vol. ۶۱, pp. ۱۲۲-۱۲۸, ۲۰۱۶.
- [۱۱] H. Koochi, K. Kiani, User based Collaborative Filtering using fuzzy C-means, *Measurement*, Vol. ۹۱, pp. ۱۳۴-۱۳۹, ۲۰۱۶.
- [۱۲] S.H.S. Chee, J. Han, K. Wang, RecTree: An Efficient Collaborative Filtering Method, *International Conference on Data Warehousing and Knowledge Discovery, DaWaK*

