

تحلیل الگوریتمهای پالایش مشارکتی مبتنی بر کالا در ارائه خدمات به شهروند الکترونیک

محمددرزی^۱، حبیب اله اصغری^۲، آندره شلتز (Szelc Andrzej Stanislaw)^۳

^۱عضو شورای علمی گروه ITBM پژوهشکده ICT جهاددانشگاهی

Modarzi@yahoo.com

^۲عضو هیات علمی و رییس مرکز رشد فناوری اطلاعات و ارتباطات جهاددانشگاهی

Asghari@itincubator.com

^۳عضو هیات علمی و معاون بین الملل دانشگاه UITM

Aszelc@wsiz.rzeszow.pl

۱- مقدمه:

حجم اطلاعات در دنیای امروز بسیار سریعتر از قدرت و توان ما در پردازش اطلاعات، روبه افزایش است. همه ما هر ساله با تولید تعداد جدیدی کتاب، مقالات و نشریات و ویژهنامه‌های کنفرانسی، افزایش روزافزون اطلاعات را احساس می‌کنیم. تکنولوژی، محدودیتهای چاپ و توزیع اطلاعات را کاهش داده است. هم اکنون زمان آن رسیده که تکنولوژیهای ایجاد شوند که بتوانند درغربال کردن اطلاعات موجود برای پیدا کردن آنچه که برای ما با ارزش تر است، به ما کمک کنند.

یکی از این تکنولوژی‌ها، سیستمهای پیشنهاددهنده است. تکنولوژی‌های متعددی در این سیستمها مورد استفاده قرار می‌گیرد که الگوریتمهای پالایش مشارکتی یکی از آنهاست [14,16,19,27]. الگوریتم پالایش مشارکتی با تشکیل پایگاه داده اولیتهای برای کالاها توسط کاربران، اجرا می‌شود. این تکنولوژی هم در حوزه تحقیق و هم در پیاده سازی وهم در کاربردهای پالایش اطلاعات و هم تجارت الکترونیکی بسیار موفق بوده است. هرچند سوالات تحقیقی مهمی در مقابله با دو چالش اساسی سیستمهای پیشنهاد دهنده پالایش مشارکتی باقیمانده است.

بحث و سوال اول، بهبود مقیاس پذیری الگوریتمهای پالایش مشارکتی است. این الگوریتمها قادرند در میان دهها هزار همسایه احتمالی بلادرنج جستجو کنند، اما سیستمهای مدرن، نیازمند جستجو در میان دهها میلیون از همسایههای احتمالی هستند. علاوه براین، الگوریتمهای موجود مورد استفاده در سایتها که حجم زیادی از اطلاعات را در اختیار دارند، مشکلات اجرایی در مورد کاربران خود دارند.

بحث دوم، در مورد بهبود کیفیت پیشنهادات برای کاربران است. کاربران به پیشنهاداتی نیاز دارند که بتوانند برای کمک در پیدا کردن کالاهایی که دوست دارند، به آنها اعتماد کنند. کاربران با رد استفاده از سیستمهای پیشنهاد دهنده که پیشنهادات آن متناسب با سلیقه آنها نیست، عملاً بی‌اعتنایی خود را به اینگونه سیستمها نشان می‌دهند.

از برخی زوایا، این دو بحث با هم در تعارض است، به طوری که الگوریتمی که زمان کمی برای جستجوی همسایهها صرف می‌کند، مقیاس پذیرتر خواهد بود ولی کیفیتش پایین تر. به همین علت، مهم است که این دو مساله به طور همزمان مورد بحث قرار گیرند و راه حلهای پیدا شده برای هر دو، مفید و عملی باشد.

در این مقاله، این مشکلات، از طریق به کار گیری الگوریتمهای مبتنی بر کالا^۱ بررسی می‌شود. معضل متداول در الگوریتمهای پالایش مشارکتی، جستجو برای همسایهها در بین جمعیت زیاد کاربران همسایه است [12]. الگوریتمهای مبتنی بر کالا از این معضل جلوگیری می‌کنند، به این صورت که ابتدا ارتباط بین کالاها بدست می‌آید تا ارتباط بین کاربران.

پیشنهادات برای کاربران از طریق پیدا کردن کالاهایی که مشابه دیگر کالاهایی هستند که کاربر به آنها علاقه داشته، محاسبه می‌شود. چون ارتباطات بین کالاها نسبتاً ایستا است، الگوریتمهای مبتنی بر کالا می‌توانند کیفیت مشابهی با الگوریتمهای مبتنی بر کاربر^۲ ولی با محاسبات برخط^۳ کمتر، تولید کنند.

۱-۱- تاریخچه:

در این قسمت، ما پیشینه کوتاهی از ادبیات تحقیقی مرتبط با پالایش مشارکتی، سیستمهای پیشنهاددهنده، داده کاوی و شخصی سازی را ارائه می‌کنیم. Tapestry [10] یکی از قدیمی ترین سیستمهای پیشنهاددهنده است که مبتنی بر پالایش مشارکتی بود. این سیستم مبتنی بر ایده های آشکار افراد یک گروه مانند یک گروه کاری در یک اداره بود. به هر حال یک سیستم پیشنهاددهنده برای یک جامعه بزرگ نمی‌تواند وابسته به فرد یا افرادی باشد که آنها دیگران را بشناسند. پس از این سیستمهای پیشنهاددهنده خودکار متعددی توسعه یافتند. سیستم پژوهشی [16,19] GroupLens یک راه حلی را برای

فیلمها و اخبار یوس نت^۴ از طریق پالایش اطلاعات ارایه کرد. Ringo [27] و پیشنهاددهنده Video [14] سیستمهای مبتنی بر پست الکترونیک و وب هستند که به ترتیب برای موزیک و فیلم پیشنهاد ارایه می کنند.

همچنین تکنولوژیهای دیگری مانند شبکه های بیژین^۵، خوشه بندی^۶ و گراف هرتینگ^۷ در سیستمهای پیشنهاددهنده به کار گرفته شده است. شبکه های بیژین یک مدل را بر مبنای یک مجموعه آموزشی^۸ با یک درخت تصمیم ایجاد می کند که هر گره و یا ل بیان کننده اطلاعات کاربر است. این مدل می تواند به صورت غیر بر خط^۹ ساخته شود. نتیجه مدل بسیار کوچک، سریع و ضرورتاً به اندازه روشهای نزدیک ترین همسایه^{۱۰} دقیق است [6]. تکنیک های خوشه بندی با تعیین کردن گروههای کاربران که مشخص می شود سلاقی مشترک دارند، عمل می کند. از زمانی که خوشه ها ایجاد می شوند، پیش بینی ها برای یک فرد می تواند از طریق گرفتن میانگین ایده های دیگر کاربران حاضر در آن خوشه ایجاد شود. میزان پیشنهادات شخصی^{۱۱} حاصل از تکنیک های خوشه بندی معمولاً از دیگر روشها کمتر است و در برخی از موارد میزان دقت خوشه ها از الگوریتمهای نزدیک ترین همسایه پایین تر است [6].

هرتینگ یک تکنیک مبتنی بر گراف است که در آن گره ها نشاندهنده کاربران است و یالهای بین گره ها بیانگر میزان شباهت بین دو کاربر است [1]. در یک مطالعه با استفاده از داده های غیر واقعی ۱۲، هورتینگ پیشنهادات بهتری را نسبت به الگوریتم نزدیک ترین همسایه ارایه کرد [1]. اسکافرو همکارانش^{۱۲} [26] طبقه بندی و مثالهایی را از سیستمهای پیشنهاددهنده که در تجارت الکترونیک به کار می روند و همچنین اینکه چطور آنها فضای شخصی سازی یک به یک^{۱۴} را فراهم می کنند و در همان زمان وفاداری مشتری^{۱۵} را بدست می آورند، ارایه کردند. اگر چه این سیستمها در گذشته موفق بوده اند ولی استفاده گسترده از آنها کاستیهایی را نشان داده است. این کاستیها مشکلاتی مانند پراکندگی در مجموعه داده هاست^{۱۶} و یا مشکلاتی که از ابعاد بالای داده ها ناشی می شود. مشکل پراکندگی در مقالات رفرنس با شماره های ۱۱ و ۲۳ مورد بحث محققان قرار گرفته و مشکلات ناشی از ابعاد بالای داده ها و تکنیک های کاهش بعد در مقالات رفرنس با شماره های ۴ و ۲۴ مورد بحث قرار گرفته است.

۱-۲- آنچه در این مقاله ارایه می شود:

۱- تشریح الگوریتم پیش بینی مبتنی بر کالا و تعیین راههای متفاوت برای اجرای زیر وظیفه^{۱۷} آنها
۲- ارایه یک مدل پیش محاسبه تشابه کالا برای افزایش مقیاس پذیری برخط پیشنهادات مبتنی بر کالا.
در بخش بعدی مقاله، پیشینه کوتاهی از الگوریتمهای پالایش مشارکتی ارایه می شود. در ابتدا فرایند پالایش مشارکتی تشریح شده و سپس دنگرش مبتنی بر حافظه و مبتنی بر مدل مورد بحث قرار می گیرد. در ادامه برخی از چالشهای پیش روی نگرش مبتنی بر حافظه ارایه می شود. در بخش سوم مقاله، نگرش مبتنی بر کالا ارایه شده و زیروظیفه های متفاوت الگوریتم به صورت کامل تشریح می شود. در بخش چهارم نیز نتیجه گیری مقاله و برخی از تجربیات استفاده از این الگوریتم ارایه می شود.

۲- سیستم های پیشنهاد دهنده مبتنی بر الگوریتم پالایش مشارکتی

سیستم های پیشنهاد دهنده، تکنیکهای تحلیل داده را جهت کمک به کاربران برای پیدا کردن کالاهایی که در سایتهای تجاری دوست دارند بخرند، به کار می گیرند. این سیستمها از طریق ارایه نمره احتمالی پیش بینی شده^{۱۸} یا فهرستی از Top-N کالای پیشنهادی به کاربران فعال، آنان را در خرید کالاهای مورد علاقه شان کمک می کنند.

ارائه پیشنهاد در رابطه با خرید کالا از طرق مختلف می تواند ایجاد شود. پیشنهادات می توانند بر اساس ویژگیهای جمعیتی کاربران، کالاهای با فروش عالی یا عادات گذشته کاربران در خرید کالا باشند. پالایش مشارکتی در حال حاضر موفق ترین تکنیک برای پیشنهاد است [19, 27]. ایده اصلی الگوریتم های مبتنی بر پالایش مشارکتی، فراهم نمودن پیش بینی یا پیشنهادات کالایی بر اساس نظرات دیگر کاربران همفکر است. نظرات کاربران می تواند به طور صریح و آشکارا از کاربر بدست آید یا به طور ضمنی از طریق بعضی ابزارها استخراج گردد.

۲-۱- نگاه کلی بر فرایند پالایش مشارکتی

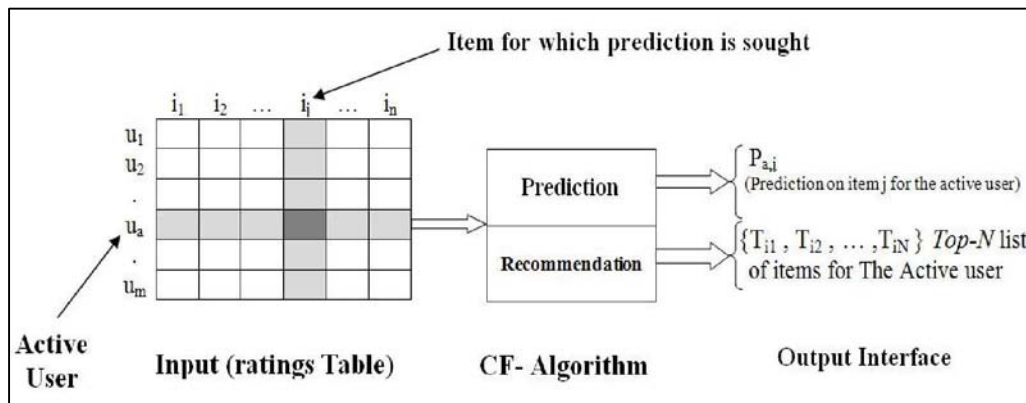
هدف الگوریتم پالایش مشارکتی، پیشنهاد کالاهای جدید یا پیش بینی مطلوبیت یک کالای معین برای یک مشتری خاص بر اساس تمایلات قبلی مشتری و نظرات دیگر مشتریان هم عقیده است. به طور کلی درسناوبی الگوریتم پالایش مشارکتی، لیستی از m کاربر به شکل $u = \{u_1, u_2, \dots, u_m\}$ و لیستی از n کالا مثل $I = \{i_1, i_2, \dots, i_n\}$ داریم. هر مشتری u_i لیستی از کالاهای I_{ui} دارد که مشتری نظراتش را در مورد آنها بیان کرده است. نظرات می توانند به طور آشکار و از طریق نمره دهی مشتری به کالاها که معمولاً با مقیاس شمارشی است، اخذ شوند و یا به طور ضمنی و بر گرفته از اطلاعات مربوط به خریدهای مشتری و یا بر طبق تحلیل ثبتهای زمانبندی شده وقایع^{۱۹} و یا کاوش ابرپیوندهای وبی^{۲۰} و ... بدست آیند [16, 28]. باید توجه داشت که I_{ui} زیر مجموعه I است و I_{ui} می تواند تهی هم باشد. در اینجا یک کاربر متمایز u_n که عضو u است، کاربر فعال نامیده می شود. وظیفه الگوریتم پالایش مشارکتی به طور مشخص، پیدا کردن کالا، متناسب با سلیقه کاربر فعال به دو شکل زیر است:

پیش بینی^{۲۱}:

پیش بینی یک مقدار عددی است مثل $P_{a,i}$ که شباهت کالای i که عضو I_{ua} نیست را برای مشتری فعال u_a بیان می کند. این مقدار پیش بینی شده با همان مقیاسی (به عنوان مثال از ۱ تا ۵) که نظرات از کاربر فعال (u_a) اخذ می شود، بیان می شود.

پیشنهاد^{۲۲}:

پیشنهاد لیستی از N کالا است؛ به عبارت بهتر I_r زیر مجموعه I است، که مشتری فعال آنها را بیشتر دوست خواهد داشت. قابل توجه آنکه که لیست پیشنهادی باید بر اساس کالاهایی باشد که قبلاً توسط مشتری فعال خریداری نشده باشد یعنی $I_r \cap I_{ua} = \Phi$. این شکل از الگوریتم های پالایش مشارکتی به عنوان پیشنهاد Top-N معروف است.



شکل 1: فرایند پالایش مشارکتی [30]

شکل ۱ شمایی از فرایند پالایش مشارکتی را نمایش می دهد. الگوریتمهای پالایش مشارکتی، داده کالا- مشتری ورودی را به عنوان ماتریس $m \times n$ رتبه A ، نمایش می دهند. هر عنصر $a_{i,j}$ در A ، نمره اولویت (رتبه بندی) کاربر i ام روی کالا j ام را نشان می دهد. رتبه بندی افراد در یک مقیاس عددی مشخص می شود و مقدار ۰، نشان دهنده این است که کاربر هنوز رتبه بندی روی کالا نداشته است. الگوریتمهای پالایش مشارکتی متعددی تاکنون معرفی شده اند که می توان آنها را به دو دسته اصلی مبتنی بر حافظه و مبتنی بر مدل تقسیم کرد [6]. در این قسمت تحلیلی از این الگوریتمها ارائه می شود.

الگوریتمهای پالایش مشارکتی مبتنی بر حافظه:

الگوریتمهای مبتنی بر حافظه از تمامی پایگاه داده کالا- مشتری برای ایجاد و آرایه یک پیش بینی استفاده می کند. این قبیل سیستمها از تکنیکهای آماری برای پیدا کردن مجموعه ای از کاربران که همسایه ها نامیده می شوند، استفاده می کنند. این کاربران در تاریخچه رفتاری خود رفتارهای مشابهی را با کاربر هدف^{۲۳} داشته اند؛ به طور مثال این همسایه ها نظرات مشابه با مشتری هدف را در خصوص کالاهای مختلف آرایه کرده و با تمایل خود را برای خرید مجموعه ای از کالای مشابه با خرید مشتری هدف نشان داده اند.

زمانی که یک همسایگی از مشتریان شکل می گیرد، این سیستمها از الگوریتمهای متفاوتی برای ترکیب سلیقه های همسایگان استفاده می کند تا پیش بینی یا پیشنهاد Top-N را به مشتری فعال آرایه کنند. این تکنیک ها همچنین به عنوان نزدیک ترین همسایه یا پالایش مشارکتی مبتنی بر کاربر معروف هستند و به طور گسترده در عمل مورد استفاده قرار می گیرند.

الگوریتمهای پالایش مشارکتی مبتنی بر مدل:

الگوریتمهای پالایش مشارکتی مبتنی بر مدل برای پیشنهاد کالا ابتدا یک مدلی را از نظرات مشتریان^{۲۴} می سازند. الگوریتمهای حاضر در این طبقه بندی از یک روش احتمالی استفاده می کنند و با فرض وجود رتبه بندیهای مشتری مورد نظر درباره سایر کالاها، از فرایند پالایش مشارکتی برای پیش بینی ارزش مورد انتظار مشتری استفاده می شود. فرایند ساخت مدل بوسیله الگوریتمهای متفاوت یادگیری ماشینی مانند شبکه بیزین، خوشه بندی و روشهای مبتنی بر قاعده^{۲۵} آرایه می شود. مدل شبکه بیزین یک مدل احتمالی را برای مساله پالایش مشارکتی فرمول بندی می کند [6]. مدل خوشه بندی با مساله پالایش مشارکتی به مثابه مشکل کلاس بندی^{۲۶} [2,6,29] رفتار می کند. روش مبتنی بر قاعده برای پیدا کردن وابستگی کالاهایی که با هم خریداری شده اند در نهایت تولید و آرایه پیشنهاد کالا بر اساس این وابستگیها از الگوریتمهای کشف قانون وابستگی استفاده می کنند [25].

۲-۲- چالش های الگوریتمهای پالایش مشارکتی مبتنی بر کاربر

سیستم های پالایش مبتنی بر کاربر، در گذشته بسیار موفق بودند، اما استفاده وسیع از آنها برخی از چالش های بالقوه را در مورد آنها آشکار کرده است که عبارتند از :

۲-۲-۱- پراکندگی^{۲۷}:

در عمل بسیاری از سیستم های پیشنهاد دهنده تجاری، برای ارزیابی مجموعه بزرگی از کالاها مورد استفاده قرار می گیرند. (به عنوان مثال Amazon.com پیشنهاد کتاب می دهد.) در این سیستم ها، حتی فعال ترین مشتریان، ممکن است کمتر از ۱٪ کالاها را خرید کرده باشند. (۱٪ از ۲ میلیون کتاب معادل ۲۰۰۰۰ کتاب می شود.) بنابراین سیستم پیشنهاددهنده بر اساس الگوریتم های نزدیکترین همسایه ممکن است قادر به ایجاد هیچ انتخابی برای کاربر خاص نباشد و در نتیجه دقت پیشنهادات پایین خواهد بود.

۲-۲-۲- مقیاس پذیری^{۲۸}:

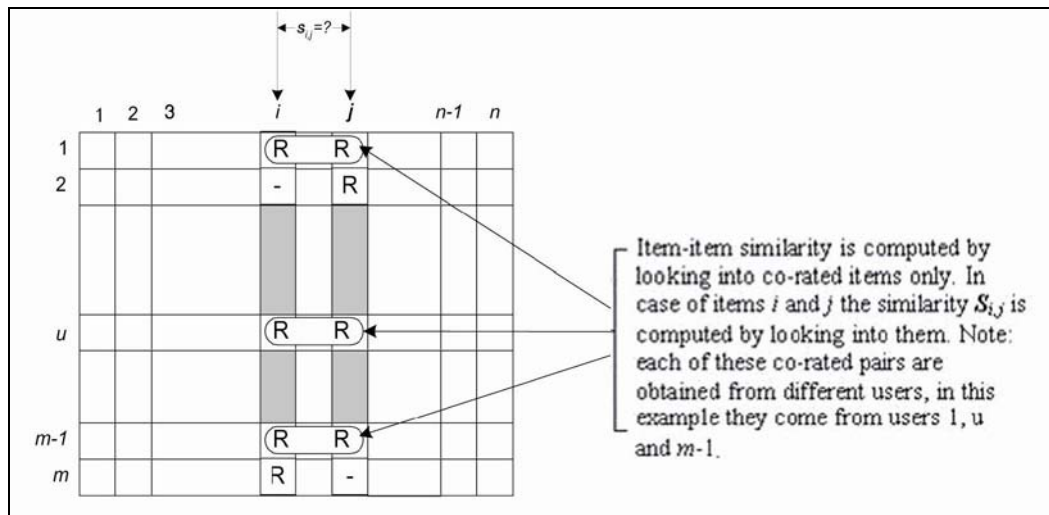
الگوریتمهای نزدیکترین همسایه، نیازمند محاسباتی هستند که با افزایش تعداد مشتریان و هم تعداد کالاها، این محاسبات افزایش می یابند. حال با داشتن میلیونها مشتری و کالا، سیستم های معمول پیشنهاد دهنده مبتنی بر وب، از مشکل جدی مقیاس پذیری رنج می برند. ضعف الگوریتم نزدیکترین همسایه برای پایگاه داده های پراکنده و بزرگ، محققان را بر آن داشت تا به سمت کشف الگوریتمهای جایگزین برای سیستمهای پیشنهاددهنده، حرکت کنند.

روش اول برای گریز از مشکل پراکندگی، استفاده از عملهای پالایش نیمه هوشمند در سیستم است [11,23]. این عملها با استفاده از ویژگی های ترکیب^{۲۹}، هر کالا را نرخ بندی و ارزیابی می کنند. از طریق ایجاد یک مجموعه رتبه بندی شده مترکم، این عملها به مشکل پوشش و بهبود کیفیت کمک می کنند. روش دیگر به کارگیری یک راهکار الگوریتمیک و استفاده از شاخص گذاری معنایی پنهان (LSI) برای بدست آوردن شباهت بین مشتریان و کالاها در فضای کاهش ابعاد است [24,25].

اما در اینجا به تکنیک دیگری که همانا روش مبتنی بر مدل است و مخصوصاً در بحث مقیاس پذیری بسیار موثر است، نگاهی می اندازیم. ایده مهم در این تکنیک، تحلیل ماتریس کالا-مشتری برای تعیین روابط بین کالاهای مختلف و سپس استفاده از این ارتباطها برای محاسبه رتبه پیش بینی یک جفت کالا-مشتری مفروض است. ایده این روش این است که یک مشتری علاقمند به خرید کالاهایی بوده که مشابه کالاهایی است که قبلاً دوست داشته و از کالاهایی دوری می کند که مشابه کالاهایی بوده که کاربر قبلاً دوست نداشته و از آنها اجتناب می ورزیده است. این تکنیک ها نیازمند تعیین همسایگی مشتریان مشابه در زمان درخواست پیشنهاد نیستند و در نتیجه آنها متمایل به تولید و ارائه پیشنهادات با سرعت بالاتر هستند. تعدادی متنوعی از روشهای مختلف اعم از احتمالی [6] تاروشهای مبتنی بر تعیین همبستگی بین کالا - کالا برای محاسبه روابط بین کالاها پیشنهاد شده اند [13,15].

۳- الگوریتم پالایش مشارکتی مبتنی بر کالا

در این قسمت یک نوع از الگوریتمهای پیشنهاد مبتنی بر کالا را برای تولید پیش بینی برای کاربر بررسی می کنیم. برخلاف الگوریتم پالایش مشارکتی مبتنی بر کاربر که در قسمت دوم این مقاله بحث شد، روش مبتنی بر کالا به مجموعه ای از کالاهایی که کاربر فعال رتبه بندی کرده توجه دارد و محاسبه می کند چه مشابهت هایی با کالای هدف i دارد و سپس k کالای بیشتر مشابه $\{i_1, i_2, \dots, i_k\}$ را انتخاب می کند. در همان زمان مشابهت های $\{s_{i_1}, s_{i_2}, \dots, s_{i_k}\}$ محاسبه می شوند. زمانی که مشابه ترین کالاها پیدا می شوند، با بدست آوردن میانگین وزنی از رتبه بندی های کاربر هدف روی این کالاهای مشابه، پیش بینی به دست می آید. در این قسمت دو مرحله از این الگوریتم شامل "محاسبه مشابهت" و "تولید پیش بینی" با جزئیات بیان می شود.



شکل 2: جداسازی کالاهای جفت رتبه بندی شده و محاسبه مشابهت [30]

۳-۱-۱- محاسبه تشابه کالا^{۳۰}:

یک گام مهم در الگوریتم های پالایش مشارکتی مبتنی بر کالا محاسبه تشابه بین کالاها و سپس انتخاب شبیه ترین کالاها است. ایده اساسی در محاسبه مشابهت بین دو کالا i و j ز در ابتدا، مجزا کردن مشتریانی است که هر دو این کالاها را رتبه بندی کرده اند و سپس به کاربردن یک تکنیک محاسبه مشابهتی برای تعیین مشابهت s_{ij} است. شکل ۲ این فرایند را توضیح داده است. در این شکل سطرهای ماتریس، مشتریان را نمایش می دهند و ستون ها هم کالاها را به نمایش می گذارند.

راههای مختلفی برای محاسبه مشابهت بین کالاها وجود دارد. در اینجا سه روش بیان می شود:

- مشابهت براساس کسینوس
- مشابهت براساس همبستگی
- مشابهت کسینوسی تنظیم شده.

۳-۱-۱-۱- مشابهت بر اساس کسینوس^{۳۱}:

در این مورد، دو کالا به عنوان دو بردار در فضای کاربری m بعدی در نظر گرفته می شوند. مشابهت بین آنها بوسیله محاسبه کسینوس زاویه بین دو بردار، اندازه گیری می شود. در ماتریس رتبه بندی $m \times n$ ، شکل ۲، مشابهت بین کالاهای i و j بوسیله $sim(i, j)$ مشخص شده و به صورت زیر است:

$$sim(i, j) = \cos(\vec{i}, \vec{j}) = \frac{\vec{i} \cdot \vec{j}}{\|\vec{i}\|_2 * \|\vec{j}\|_2} \quad (1)$$

که علامت "·" ضرب نقطه ای بین دو بردار را نمایش می دهد.

۳-۱-۲- مشابهت بر اساس همبستگی:

در این حالت، مشابهت بین دو کالای i و j از طریق محاسبه ضریب همبستگی پیرسون^{۳۲} به صورت $corr_{ij}$ ، اندازه گیری می شود. برای ایجاد دقت محاسباتی همبستگی، ابتدا باید موارد با هم رتبه بندی شده را جدا کرد. (مثلا مواردی که مشتری هر دو کالا i و j را رتبه بندی کرده است). مجموعه ای از مشتریان را که هر دو کالا i و j را رتبه بندی کرده اند U در نظر می گیریم و مشابهت همبستگی به صورت زیر بدست می آید:

$$sim(i, j) = \frac{\sum_{u \in U} (R_{u,i} - \bar{R}_i)(R_{u,j} - \bar{R}_j)}{\sqrt{\sum_{u \in U} (R_{u,i} - \bar{R}_i)^2} \sqrt{\sum_{u \in U} (R_{u,j} - \bar{R}_j)^2}} \quad (2)$$

در اینجا $R_{u,i}$ رتبه بندی کاربر u روی کالا i و \bar{R}_i میانگین رتبه بندی کالای i ام است.

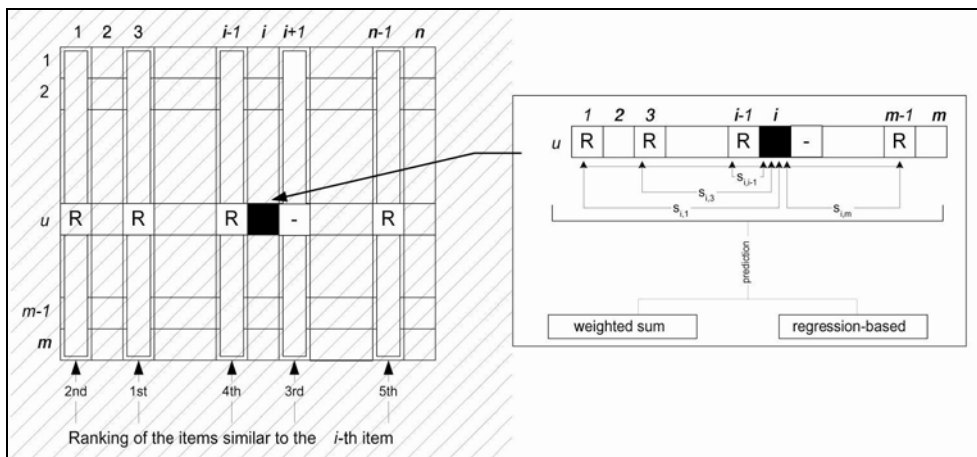
۳-۱-۳- مشابهت کسینوسی تنظیم شده^{۳۳}:

یک تفاوت اساسی بین محاسبه مشابهت در پالایش مشارکتی مبتنی بر کاربر با پالایش مشارکتی مبتنی بر کالا در این است که مشابهت در پالایش مشارکتی مبتنی بر کاربر، در طول سطرهای ماتریس محاسبه انجام می شود اما در مورد پالایش مشارکتی مبتنی بر کالا در طول ستون. محاسبه تشابه با استفاده از روش کسینوسی در موارد مبتنی بر کالا یک اشکال مهم دارد و آنهم این است که تفاوت در مقیاس رتبه بندی شده بین کاربران مختلف لحاظ نمی شود. مشابهت کسینوسی تنظیم شده این مانع را بوسیله کم کردن مقدار میانگین کاربر متناظر با هر جفت کالای رتبه بندی شده، متعادل می کند.

فرمول مشابهت بین کالای i و j مورد استفاده به شکل زیر است :

$$sim(i, j) = \frac{\sum_{u \in U} (R_{u,i} - \bar{R}_u)(R_{u,j} - \bar{R}_u)}{\sqrt{\sum_{u \in U} (R_{u,i} - \bar{R}_u)^2} \sqrt{\sum_{u \in U} (R_{u,j} - \bar{R}_u)^2}} \quad (3)$$

در اینجا \bar{R}_u رتبه بندی کاربر u است.



شکل 3: فرایند تولید پیش بینی برای ۵ همسایه در پالایش مشارکتی مبتنی بر کالا [30]

۳-۲- محاسبه پیش بینی^{۳۴}:

مهمترین مرحله در یک سیستم پالایش مشارکتی، تولید رابط خروجی^{۳۵} مربوط به پیش بینی است. وقتی که مجموعه ای از شبیه ترین کالاها براساس ابزارهای تعیین مشابهت مجزا می شوند، مرحله بعدی مشاهده رتبه بندی های مشتریان هدف و استفاده از یک تکنیک برای به دست آوردن پیش بینی ها است. در اینجا دو تکنیک بررسی می شود :

۳-۲-۱- مجموع وزنی^{۳۶}:

همانطور که از نام تکنیک مشخص است این روش پیش بینی روی یک کالای i را برای کاربر u از طریق محاسبه مجموع رتبه های داده شده بوسیله مشتری روی کالاهای مشابه i ، محاسبه می کند. هر رتبه بندی بوسیله مشابهت $S_{i,j}$ متناظر بین کالاهای i و j وزن دهی می شود. به صورت تفصیلی این نکات در شکل ۳ مشخص شده اند و پیش بینی $p_{u,i}$ به صورت زیر تعیین می شود:

$$P_{u,i} = \frac{\sum_{\text{all similar items } , N} (S_{i,N} * R_{u,N})}{\sum_{\text{all similar items } , N} |S_{i,N}|} \quad (4)$$

اساساً این روش سعی دارد تأکید ارزیابی کاربر فعال روی کالاهای مشابه را بدست آورد. مجموع وزنی بوسیله مجموع شرایط مشابهت برای اطمینان از اینکه پیش بینی مطابق محدوده از قبل تعریف شده باشد، اندازه گیری می شود.

۳-۲-۲- رگرسیون:

این روش مشابه متد میانگین وزنی است اما به جای استفاده مستقیم رتبه بندی های کالاها از تقریبی از رتبه بندی ها براساس مدل رگرسیون استفاده می کند. در عمل، مشابهت های محاسبه شده با استفاده از کسینوس یا روش همبستگی ممکن است گمراه کننده باشد؛ بدین صورت که در حالتی که دو بردار رتبه^{۲۷}، فاصله دار هستند (در مفهوم اقلیدسی) هنوز مشابهت های خیلی زیادی دارند. در این مورد استفاده از رتبه بندی های سطری به اصطلاح "کالای مشابه"، ممکن است باعث پیش بینی ضعیفی گردد. ایده اساسی، استفاده از فرمول مشابه فرمول مجموع وزنی است، اما به جای استفاده مقادیر رتبه بندی N سطری کالا مشابه این مدل از مقدار تقریبی $R_{u,N}$ خودش، براساس مدل رگر $R'_{u,N}$ خطی استفاده می کند. اگر بردارهای کالای هدف i و کالای مشابه N بوسیله R_i و R_N مشخص شوند، مدل رگرسیون خطی می تواند به صورت زیر بیان شود:

$$\bar{R}'_N = \alpha \bar{R}_i + \beta + \epsilon \quad (5)$$

α و β پارامترهای مدل رگرسیون هستند که با حرکت بر روی هر دو بردارهای رتبه بندی تعیین می شوند. ϵ خطای پارامتر مدل رگرسیون است.

۴- نتیجه گیری:

سیستمهای پیشنهاددهنده فناوری قدرتمند و جدیدی هستند که از طریق پایگاه داده مشتریان، باعث ایجاد ارزش افزوده در کسب و کارها می گردد. این سیستم ها مشتریان را در پیدا کردن کالاهایی که می خواهند خریداری نمایند، کمک می نمایند. سیستمهای پیشنهاددهنده با توانمند ساختن مشتریان در خرید کالاهای مورد علاقه شان به مشتریان سود می رسانند و در مقابل، این سیستمها کسب و کارها را با ایجاد فروش بیشتر یاری می رسانند. سیستمهای پیشنهاددهنده به سرعت در حال تبدیل شدن به یکی از ابزارهای حیاتی سایتهای تجارت الکترونیک در محیط وب هستند. برای بهبود مساله مقیاس پذیری در سیستمهای پیشنهاددهنده نیاز به تکنولوژیهای جدید است.

در این مقاله، الگوریتمهای پالایش مشارکتی که در سیستمهای پیشنهاددهنده به کار می رود تشریح شدند. نتایج تحقیقات [30] و تجربیات شرکتی که این الگوریتم استفاده کرده اند [31]، نشان از موفقیت این الگوریتمها دارد.

مراجع:

- [1] Aggarwal, C. C., Wolf, J. L., Wu K., and Yu, P. S. (1999). *Horting Hatches an Egg: A New Graph-theoretic Approach to Collaborative Filtering*. In Proceedings of the ACM KDD'99 Conference. San Diego, CA. pp. 201-212.
- [2] Basu, C., Hirsh, H., and Cohen, W. (1998). *Recommendation as Classification: Using social and Content-based Information in Recommendation*. In Recommender System Workshop '98. pp. 11-15.
- [3] Berry, M. W., Dumais, S. T., and O'Brian, G. W. (1995). *Using Linear Algebra for Intelligent Information Retrieval*. SIAM Review, 37(4), pp. 573-595.
- [4] Billsus, D., and Pazzani, M. J. (1998). *Learning Collaborative Information Filters*. In Proceedings of ICML '98. pp. 46-53.
- [5] Brachman, R., J., Khabaza, T., Kloesgen, W., Piatetsky-Shapiro, G., and Simoudis, E. 1996. *Mining Business Databases*, Communications of the ACM, 39(11), pp. 42-48, November.
- [6] Breese, J. S., Heckerman, D., and Kadie, C. (1998). *Empirical Analysis of Predictive Algorithms for Collaborative Filtering*. In Proceedings of the 14th Conference on Uncertainty in Artificial Intelligence, pp. 43-52.
- [7] Cureton, E. E., and D'Agostino, R. B. (1983). *Factor Analysis: an Applied Approach*, Lawrence Erlbaum associates pubs, Hillsdale, NJ.
- [8] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990). *Indexing by Latent Semantic Analysis*, Journal of the American Society for Information Science, 41(6), pp. 391-407.
- [9] Fayyad, U. M., Piatetsky-Shapiro, G., Smyth, P., and Uthurusamy, R., Eds. (1996). *Advances in Knowledge Discovery and Data Mining*, AAAI press/MIT press.
- [10] Goldberg, D., Nichols, D., Oki, B. M., and Terry, D. (1992). *Using Collaborative Filtering to Weave an Information Tapestry*, Communications of the ACM, December.
- [11] Good, N., Schafer, B., Konstan, J., Borchers, A., Sarwar, B., Herlocker, J., and Riedl, J. (1999). *Combining Collaborative Filtering With Personal Agents for Better Recommendations*, In Proceedings of the AAAI-'99 conference, pp 439-446.
- [12] Herlocker, J., Konstan, J., Borchers, A., and Riedl, J. (1999). *An Algorithmic Framework for Performing Collaborative Filtering*. In Proceedings of ACM SIGIR'99. ACM press.
- [13] Herlocker, J. (2000). *Understanding and Improving Automated Collaborative Filtering Systems*, Ph.D. Thesis, Computer Science Dept., University of Minnesota.
- [14] Hill, W., Stead, L., Rosenstein, M., and Furnas, G. (1995). *Recommending and Evaluating Choices in a Virtual Community of Use*. In Proceedings of CHI '95.
- [15] Karypis, G. (2000). *Evaluation of Item-Based Top-N Recommendation Algorithms*, Technical Report CS-TR-00-46, Computer Science Dept., University of Minnesota.
- [16] Konstan, J., Miller, B., Maltz, D., Herlocker, J., Gordon, L., and Riedl, J. (1997). *GroupLens: Applying Collaborative Filtering to Usenet News*, Communications of the ACM, 40(3), pp. 77-87.
- [17] Ling, C. X., and Li C. (1998). *Data Mining for Direct Marketing: Problems and Solutions*, In Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining, pp. 73-79.

- [18] Peppers, D., and Rogers, M. (1997). *The One to One Future: Building Relationships One Customer at a Time*, Bantam Doubleday Dell Publishing.
- [19] Resnick, P., Iacovou, N., Suchak, M., Bergstrom, P., and Riedl, J. (1994). *GroupLens: An Open Architecture for Collaborative Filtering of Netnews*. In Proceedings of CSCW '94, Chapel Hill, NC.
- [20] Resnick, P., and Varian, H. R. (1997). *Recommender Systems*, Communications of the ACM, 40(3).
- [21] Reichheld, F. R., and Sasser Jr., W. (1990). *Zero Defections: Quality Comes to Services*, Harvard Business School Review, 1990(5): pp. 105-111.
- [22] Reichheld, F. R. (1993). *Loyalty-Based Management*, Harvard Business School Review, 1993(2): pp. 64-73.
- [23] Sarwar, B. M., Konstan, J. A., Borchers, A., Herlocker, J., Miller, B., and Riedl, J. (1998). *Using Filtering Agents to Improve Prediction Quality in the GroupLens Research Collaborative Filtering System*, In Proceedings of CSCW '98, Seattle, WA.
- [24] Sarwar, B. M., Karypis, G., Konstan, J. A., and Riedl, J. (2000). *Application of Dimensionality Reduction in Recommender System—A Case Study*, In ACM WebKDD 2000 Workshop.
- [25] Sarwar, B. M., Karypis, G., Konstan, J. A., and Riedl, J. (2000). *Analysis of Recommendation Algorithms for E-Commerce*, In Proceedings of the ACM EC'00 Conference. Minneapolis, MN, pp. 158-167
- [26] Schafer, J. B., Konstan, J., and Riedl, J. (1999). *Recommender Systems in E-Commerce*. In Proceedings of ACM E-Commerce 1999 conference.
- [27] Shardanand, U., and Maes, P. (1995). *Social Information Filtering: Algorithms for Automating 'Word of Mouth'*, In Proceedings of CHI '95. Denver, CO.
- [28] Terveen, L., Hill, W., Amento, B., McDonald, D., and Creter, J. (1997). *PHOAKS: A System for Sharing Recommendations*, Communications of the ACM, 40(3), pp. 59-62.
- [29] Ungar, L. H., and Foster, D. P. (1998) *Clustering Methods for Collaborative Filtering*, In Workshop on Recommender Systems at the 15th National Conference on Artificial Intelligence.
- [30] Sarwar, B., Karypis, G., Konstan, J. and Riedl, J. (2001), *Item-based Collaborative Filtering Recommendation Algorithms*, Proceedings of the 10th international conference on World Wide Web Hong Kong, pp: 285 – 295.
- [31] Linden, G., Smith, B., and York, J. (Jan/Feb 2003), *Amazon.com recommendations: item-to-item collaborative filtering*, In IEEE Internet computing, Volume: 7, pp: 76- 80.

-
- 1 - Item-based Algorithms
 - 2 - User-based Algorithms
 - 3 - Online
 - 4 - Usenet news
 - 5 - Bayesian networks
 - 6 - Clustering
 - 7 - Horting graph
 - 8 - Training set
 - 9 - Off-line
 - 10 - Nearest neighbor
 - 11 - Personal recommendation
 - 12 - Synthetic data
 - 13 - Schafer et al.
 - 14 - One-to-one personalization
 - 15 - Customer loyalty
 - 16 - Sparsity in data set
 - 17 - Subtasks
 - 18 - predicted likeliness score
 - 19 - Analyzing timing logs
 - 20 -Web mining hyperlinks
 - 21 - Prediction
 - 22 - Recommendation
 - 23 - Target user
 - 24 - User rating
 - 25 - Rule – based methods
 - 26 - Cclassification
 - 27 - Sparsity
 - 28 - Scalability
 - 29 - Syntactic Features
 - 30 - Item Similarity Computation
 - 31 - Cosine-based similarity
 - 32 - Pearson Correlation
 - 33 - Adjusted Cosine Similarity
 - 34 - Prediction Computation
 - 35 - Output Interface

³⁶ - Weighted sum

³⁷ - Rating Vector