

ارائه یک سیستم توصیه گر برای استخراج اطلاعات مورد نیاز کاربران از کتابخانه های دیجیتالی بزرگ

رضا مولایی فرد

کارشناسی ارشد کامپیوتر نرم افزار

نام نویسنده مسئول:

رضا مولایی فرد

تاریخ دریافت: ۱۳۹۸/۱۱/۲۳

تاریخ پذیرش: ۱۳۹۹/۱/۱۹

چکیده

امروزه کتابخانه های دیجیتال یکی از بهترین راهکارها برای به دست آوردن اطلاعات مورد نیاز کاربران می باشد. این کتابخانه های دیجیتال با توجه به دسترسی آسان تر و هزینه کمتر یکی از بهترین راهکارهای کاربران برای بدست آوردن کتاب های موردنظر کاربران می باشد ولی با توجه به رشد روز افزون تعداد کتاب های مختلف در زمینه های مختلف وجود سیستمی که بتواند اطلاعات مورد نیاز افراد را از میان حجم عظیم این کتاب ها که روز به روز در حال افزایش می باشند، استخراج کند لازم و ضروری به نظر می رسد، یکی از بهترین راهکارها برای بدست آوردن کتب موردعلاقه کاربران از میان حجم عظیم داده ها، شخصی سازی این سیستم ها می باشد. یکی از بهترین سیستم های شخصی سازی استفاده از سیستم های توصیه گر می باشند، سیستم های توصیه گر یا پیشنهاد دهنده سیستم هایی هستند که با یکسری از روش های داده کاوی و روش هایی مانند اینکه کاربر در گذشته چه کتاب هایی جستجو کرده است می تواند پیشنهاد های مناسبی به کاربر ارائه دهند. در این تحقیق از یک سیستم توصیه گر مبتنی بر پالایش مشارکتی استفاده شده است که با دسته بندی کاربران مشابه در یک گروه و با توجه به سوابق این افراد و امتیازهایی که به داده ها در پروفایل خود داده اند و تشابهات بین آنها به کاربر موردنظر می تواند پیشنهاد مناسبی از بین این حجم کتب مختلف ارائه دهد که موردنظر و مورد علاقه کاربر باشد.

واژگان کلیدی: سیستم توصیه گر، استخراج اطلاعات، داده کاوی، کتابخانه دیجیتالی، فیتینگ مشارکتی.

مقدمه

با توجه به رشد روز افزون داده‌ها و کتب و همچنین با توجه به مشکلاتی که جهت تهیه کتب مختلف وجود دارد، وجود یک سیستم الکترونیکی که بتواند مجموعه‌ای از این کتب را در اختیار کاربران قرار داده و مشکلات دسترسی به کتب را حل کند لازم و ضروری به نظر می‌رسد. کتابخانه‌های دیجیتال راه حل مناسبی برای دسترسی به کتب مورد نیاز افراد می‌باشد. کتابخانه‌های دیجیتالی یکی از زمینه‌های پژوهشی است که امروزه در برخی از رشته‌ها دنبال می‌شود. کوشا در مقاله خود تعاریف مختلف را که نتیجه چهار دیدگاه مخلف در زمینه کتابخانه هستند شناسایی نموده که عبارتند از: تعاریف ارائه شده توسط متخصصان رایانه: در این دیدگاه کتابخانه دیجیتال در قالب پایگاه اطلاعاتی و نظام بازیابی اطلاعات تعریف می‌شود: تعاریف ارائه شده توسط محققان کتابداری و اطلاع‌رسانی که در آن کتابخانه دیجیتال به مثابه سازمان و خدمات است؛ تعرف ارائه شده توسط سازمان‌های مجری و سیاست‌گذار که به دو دسته تعاریف ارائه شده از طرف سازمان‌های سیاست‌گذار کتابخانه‌های دیجیتال مانند فدراسیون بین‌المللی کتابخانه‌های دیجیتال و تعاریف ارائه شده توسط سازمان‌های مجری کتابخانه‌های دیجیتال تقسیم می‌شوند؛ و تعاریف پژوهش‌مدار که در آن‌ها کتابخانه دیجیتال به منزله یک مسئله پژوهشی در نظر گرفته و تعریف می‌شوند. جوامع مختلفی در امر کتابخانه دیجیتال دخیلند و هرکدام رویدادها، مفاهیم، معانی، تمرکز و رویکرد خود را دارند. دو گروه عمده که به پژوهش در زمینه کتابخانه‌های دیجیتال می‌پردازند عبارتند از علم و مهندسی کامپیوتر و علم و اطلاعات و دانش‌شناسی. برخی از موضوعات مورد توجه در عرصه علوم کامپیوتر عبارتند از: تعامل بین مجموعه‌های متجانس، یکپارچه‌سازی داده‌ها مانند ویدیو، متن، صوت، استانداردها و پروتکل شبکه‌ای، موتورهای جستجو و عامل‌های جستجو، محدودسازی، تلخیص و یکپارچه‌سازی داده‌ها را تسهیل می‌کنند، فناوری‌های مصورسازی و تعاملی مانند امکان مرور حجم وسیعی از متون یا تصاویر. رویکردهای غالب در علم اطلاعات و دانش‌شناسی نیز عبارتند از: رویکرد سازمان، خدمات و فعالیت محور: در این رویکرد مفهوم کتابخانه گسترده می‌شود؛ رویکرد محتوا و مجموعه محور: که در آن منابع اطلاعاتی مختلف و همچنین مخازن مختلف منابع دیجیتال مورد توجه قرار می‌گیرند؛ رویکرد دسترسی به مجموعه‌ها: که در این رویکرد به مسائلی نظیر محیط‌ها و جوامع کاربردی مختلف و میزان یکپارچه‌سازی یا جداسازی منابع اطلاعاتی می‌پردازد (saracevic 1999). با وجود همه این تفاسیر یک راه حل مناسب جهت استخراج اطلاعات، استفاده از سیستم‌های توصیه‌گر می‌باشند، سیستم‌های توصیه‌گر یا پیشنهاددهنده سیستم‌هایی هستند که با یکسری از روش‌های داده‌کاوی می‌توانند پیشنهادات مناسبی را به کاربر برای استخراج اطلاعات ارائه دهند. در این تحقیق از یک سیستم توصیه‌گر مبتنی بر فیلترینگ مشارکتی استفاده می‌شود که می‌تواند اطلاعات مورد نیاز کاربر را استخراج و پیشنهادات مناسبی را به وی ارائه دهند که مورد نظر کاربر باشند.

پیشینه پژوهش

کیم و لی در سال ۲۰۰۸ در ژورنالی با عنوان بررسی ساختار درونی مطالعات آرشیوی با روش متن‌کاوی در سال‌های ۲۰۰۱ تا ۲۰۰۴، به بررسی ۴۳۲ مقاله از سال ۲۰۰۱ تا ۲۰۰۴ پرداخته و ۴۳ خوشه از اسناد با استفاده از میانگین درون‌گروهی در نرم‌افزار SPSS ایجاد کرده‌اند. سپس شبکه‌های مسیریاب این ۴۳ خوشه را ساخته و در ۷ شاخه موضوعی شامل: کتابخانه‌های دیجیتال و فناوری آرشیوسازی دیجیتال، منابع آنلاین و کمک‌های اکتشافی، آرشیوها و آرشیویست‌ها، موضوعات سیاسی و حقوقی، مسائل فنی و پیشینه‌های الکترونیک، مدیریت اطلاعات و رکوردها و پست الکترونیک و اطلاعات حرفه‌ای گروه‌بندی کرده‌اند. در نهایت این ۷ موضوع در ۳ بخش ادغام شده‌اند که شامل کتابخانه‌های دیجیتال، آرشیوها و پژوهش‌ها در عمل می‌شوند. این مطالعه تغییرات پویای موضوعات پژوهشی این رشته، از حوزه‌های موضوعی صرفاً سنتی تا ظهور حوزه‌های موضوعی پیچیده را از سال ۲۰۰۱ تا ۲۰۰۴ نشان می‌دهد. نتایج این پژوهش بیان می‌کند که حوزه‌های پژوهشی در علوم آرشیوی پتانسیل رشد بالایی داشته و همچنان توسعه خواهند یافت.

عرب در سال ۱۳۹۲ در مقاله‌ای با عنوان، ارزیابی قابلیت‌های جستجو پنج‌گانه کتابخانه دیجیتالی ایرانی، به بررسی کتابخانه‌های دیجیتالی ایرانی (نور، دید، تبیان، میراث فرهنگی، صنایع دستی و گردشگری و ارم) پرداخته است. نتایج پژوهش نشان

داده در کتابخانه‌های دیجیتالی مذکور، هیچکدام به دو ویژگی جستجو، جستجوی مجاورتی و جستجوی کوتاه سازی توجه نداشته‌اند، ولی همه کتابخانه‌ها دو مورد جستجوی ساده و جستجوی پیشرفته را دارا می‌باشند.

زائو و زنگ در سال ۲۰۱۱ در مقاله‌ای با عنوان، ترسیم نقشه علمی پژوهش‌های تولید شده چینی در حوزه کتابخانه در سال‌های ۱۹۹۴ تا ۲۰۱۰ برای دست‌یابی به پارادایم‌های فکری در حوزه کتابخانه‌های دیجیتالی، با استفاده از روش‌های متن‌کاوی و تحلیل هم‌رخدادی واژگان به بررسی ۱۲۵۰ و ۶۰۶۸ مقاله‌بازیابی شده در حوزه کتابخانه دیجیتالی به ترتیب در پایگاه ملی دانش زیربنایی چین و پایگاه سانس دایرکت پرداخت. در نهایت آنها با استفاده از فنون تحلیل شبکه‌های اجتماعی، نقشه‌های موضوعی و الگوهای پژوهشی کتابخانه‌های دیجیتال در چین را با استفاده از نرم‌افزارهای UCINET و Netdraw ترسیم کردند.

ساروش پاریک در سال ۲۰۱۳ در مقاله خود به صورت کتاب‌سنجی به مطالعه متون مجلات ایفلا مربوط به سال‌های ۲۰۰۱ تا ۲۰۱۰ پرداخته و هدف از مطالعه خود را تحلیل کتاب‌سنجی از ویژگی‌هایی مانند میزان انتشار مقاله در سال، الگوهای نویسندگی، میزان همکاری موسسات، توزیع موضوعات، الگوی استناد و... بیان کرده است.

بالاجی بابو و کریشنامورسی در سال ۲۰۱۳ در پژوهشی با هدف مطالعه‌ی تجزیه و تحلیل الگویی از اتوماسیون کتابخانه برای کشف منابع بوسیله کاوش برنامه‌های کاربردی کشف منابع به این نتیجه رسید که، رشد صنعت اتوماسیون هند پر رونق است. با این حال سازگاری نرم‌افزار کتابخانه و پیشرفت جامعه رضایت‌بخش نیستند.

نوروزی در سال ۱۳۸۹ در پژوهشی تحت عنوان، بررسی میزان رعایت معیارهای ارزیابی رابط کاربر در صفحات وب فارسی کتابخانه‌های دیجیتالی خود ساخته و خریداری شده در ایران، به ارزیابی رابط کاربر کتابخانه‌های دیجیتالی ایران پرداخته است. روش انجام پژوهش از نوع تحقیقات ارزیابانه است که از یک سیاهه واری متشکل از ۱۰ معیار اصای و ۱۱۴ مولفه فرعی برای ارزیابی و تحلیل ابعاد پژوهش استفاده شده است. یکی از معیارهای اصلی این پژوهش، جستجو می‌باشد. این پژوهش نشان می‌دهد که از همه امتیازات مربوط به معیار ده‌گانه، کتابخانه‌های دیجیتالی خودساخته توانستند در هفت معیار تصحیح خطا، کنترل کاربر، راهبری، جستجو و سادگی امتیاز بالایی را کسب کنند. در صورتی که کتابخانه‌های دیجیتالی خریداری شده، تنها در سه معیار زبان رابط کاربر، نمایش اطلاعات و انسجام وضعیت بهتری قرار دارند.

سیستم توصیه گر

سیستم‌های توصیه گر^۱ یا پیشنهاد دهنده زیر مجموعه‌ای از سامانه پالایش اطلاعات که به دنبال پیش‌بینی امتیاز یا اولوبیتی است که کاربر به یک آیتم (داده، اطلاعات، کالا و...) خواهد داد. در سال‌های اخیر سیستم‌های توصیه گر بسیار متداول شده و در حوزه‌های مختلفی مورد استفاده قرار گرفته‌اند. برخی از کاربردهای معروف این سیستم‌ها می‌توان به موارد زیر اشاره کرد: موسیقی، صفحات وب، اخبار، کتاب‌ها و مقالات، صنعت گردشگری، صنعت هتل‌داری، صنعت پزشکی، جست و جو و شبکه‌های اجتماعی. در واقع سامانه‌های توصیه گر مانند یک فیلتر عمل می‌کنند، فیلتری که فقط آنچه مطلوب کاربر است را به او نشان می‌دهد که به این عمل شخصی سازی کردن اطلاعات می‌گویند. به طور کل سامانه توصیه گر یک سامانه پشتیبان شخصی سازی است که اطلاعات را به سه عامل تعیین کننده، سفارشی سازی، علاقه مندی و سودمندی، برای کاربران ویژه بوسیله تجزیه و تحلیل سلیق آنها و محتوای آیتم‌ها می‌سنجد. سیستم توصیه گر از جمله ابزارهایی است که می‌تواند کاربران را در محیط‌های الکترونیکی به سمت یافتن اطلاعات، خدمات و آیتم‌های مورد نظرشان هدایت کند [کونستانت و همکاران، ۲۰۰۸]. سیستم‌های توصیه گر با قابلیت کشف علائق کاربران و پیش‌بینی اولوبت آنها، آیتم‌هایی که احتمال می‌رود مورد توجه کاربر باشد را از بین حجم بالای داده‌ها پالایش کرده و یا آنها را پیشنهاد آنها، در زمان او صرفه جویی می‌کند. از طرف دیگر این سیستم‌ها توانایی ذخیره و تحلیل رفتارهای گذشته کاربر، خدمات و اطلاعاتی را که مورد توجه کاربران نبوده و احتمالاً به آنها علاقه مند هستند را نیز استنتاج کرده و نتایج جالب توجهی به کاربران ارائه می‌کند. در واقع سیستم‌های توصیه گر یکی از ابزارهای اصلی غلبه بر مشکل افزونگی اطلاعات بوده و با داشتن قدرت تحلیل رفتارهای کاربر، مکمل هوشمندی

¹ Recommender system

برای مفاهیم بازایی و پالایش اطلاعات است. امروزه سیستم های توصیه گر در زمینه های متنوعی از پالایش اطلاعات موجود در وب متناسب با خواسته های کاربر تا تجارت الکترونیکی، پیشنهاد فیلم، موزیک، کتاب، مقاله و... کاربرد دارند. سیستم های توصیه گر انواع مختلفی دارند که به طور کلی به دسته های زیر تقسیم می شوند:

سیستم توصیه گر مبتنی بر محتوا

این روش بر مبنای توصیف آیتم ها با کلمات کلیدی و ایجاد پروفایل براساس اولویت های کاربر است. در واقع، این الگوریتم مواردی را براساس تشابهات با علائق کاربر در گذشته توصیه می کند، یعنی مواردی که قبلا توسط کاربر امتیازدهی شده با آیتم موجود مقایسه کرده و بهترین تطابق را توصیه می کند. این سیستم با استفاده از مدل اولویت کاربر و سابقه تعامل کاربر پروفایلی را ایجاد می کند. این سیستم همچنین پروفایلی مبتنی بر محتوا براساس بردار وزن ویژگی ها ایجاد می کند. این وزن ها دلالت بر اهمیت ویژگی ها توسط کاربر دارند و بر اساس عقیده کاربر افزایش یا کاهش می یابند (رائو کاگیتا و همکاران، ۲۰۱۵؛ لو و همکاران، ۲۰۱۵؛ وانگ و همکاران، ۲۰۱۵).

سیستم توصیه گر مبتنی بر فیلترینگ مشارکتی

مبنای فیلترینگ مشارکتی بر اولویت های رفتاری یا فعالیت های کاربران و پیش بینی علائق آنان، براساس تشابهات با کاربران دیگر است. این روش با پیش بینی خودکار فیلترینگ، علائق کاربران، اطلاعات را جمع آوری می نمایند. این روش معمولا به مشارکت کاربر، یعنی تجزیه و تحلیل پروفایل و الگوریتمی برای یافتن افراد با علائق مشابه نیاز دارد. در واقع، در این روش کاربر نظر خود را با امتیاز دهی اقلام در سیستم بیان می کند و سیستم های کاربرانی را با الگوهای امتیاز دهی یکسانی را به اشتراک می گذارند یافته و از این کاربران همفکر برای محاسبه پیش بینی استفاده می کند (پاول و همکاران، ۲۰۱۵؛ کیروش داسیلوا و همکاران، ۲۰۱۶؛ یانگ و همکاران، ۲۰۱۴؛ وی و همکاران، ۲۰۱۵).

سیستم توصیه گر مبتنی بر دانش

این سیستم ها بر اساس ادراکی که از نیاز های کاربر و ویژگی های آیتم ها پیدا کرده اند، توصیه ارائه می دهند. به عبارتی در این گونه از سیستم های توصیه گر مواد اولیه مورد استفاده برای تولید لیستی از توصیه ها، دانش سیستم در مورد کاربر و آیتم است [بورکه، ۲۰۰۰]. سیستم های مبتنی بر دانش از متد های مختلفی که برای تحلیل دانش، قابل استفاده هستند بهره میبرند که متد های رایج در الگوریتم های ژنتیک، فازی، شبکه های عصبی و... از جمله آنها. همچنین در این گونه سیستم ها از درخت ای تصمیم، استدلال های مبتنی بر شاهد و... نیز می توان استفاده کرد. یکی از رایج ترین متدهای مبتنی بر دانش در سیستم های توصیه گر، روش استدلال مبتنی بر نمونه است.

سیستم های توصیه گر ترکیبی

در سیستم های توصیه گر ترکیبی برای رسیدن به بالاترین کارایی، بر اساس یک استراتژی معین تکنیک های مختلف با یکدیگر ترکیب می شوند [بورکه، ۲۰۰۲]. برای مثال دو تکنیک پالایش مشارکتی و متنی بر دانش در صورتی که با یکدیگر ترکیب شوند، نتیجه سیستمی خواهد بود که به واسطه جزء مبتنی بر دانش، می تواند مشکل شروع آهسته تکنیک پالایش مشارکتی را جبران کند و از طرفی با وجود جزء پالایش مشارکتی و قدرت بالای آن در یافتن اولویت های مشابه کاربران می تواند توصیه ها را تولید کند که هیچ سیستم مبتنی بر دانش، توانایی توصیه به آن را بر اساس دانش ندارد.

داده کاوی

مرکز تحقیقات آمریکا و اداره پاسخگویی سازمان ها داده کاوی را مستلزم استفاده از ابزارهای پیشرفته برای تحلیل و کشف روابط و استخراج الگوهای ارزشمند داده ها به منظور دسترسی به قوانین جدید معنا دار می داند (مرادی و قاسمی؛ ۱۳۹۷). کشف دانش درون داده ها در عصر اطلاعات یکی از هیجان انگیزترین و کلیدی ترین مفاهیمی است که روز به روز به اهمیت آن افزوده می شود (رحمانی و حاجی زین العابدینی؛ ۱۳۹۵). با توجه به تعاریف و تفاسیر مطرح شده از دیدگاه های مختلف، می توان دو جزء اساسی را در داده کاوی مشخص نمود، اولی کشف الگوهای پنهان در داده ها می باشد و دیگری استفاده از این الگوها برای پیش بینی نتایج در آینده است (مرادی و قاسمی؛ ۱۳۹۷). در سال های اخیر با پیشرفت سریع و گسترده شبکه ها و فناوری پایگاه های اطلاعاتی و کتابخانه ها هم با حجم بسیار زیادی از اطلاعات مواجهه بوده اند. با گسترش کارایی کتابخانه های دیجیتالی و انبوه منابع اطلاعاتی ناهمگن در حال رشد که جان نشرت، معماری، غنای اطلاعاتی و فقر دانش را مطرح می نماید (پناه؛ ۱۳۹۶) در آن فناوری داده کاوی با دو اصل مهم پیش بینی نتایج و کشف دانش در اثرسنجی وب جهان گستر در این گونه کتابخانه ها و مهم تر در جذب کاربران خاص با منابع اطلاعاتی متنوع ارزش کاربردی بسیاری پیدا کرده است.

وب کاوی

عبارت وب کاوی مترادف با یکی از عبارت های استخراج دانش، برداشت اطلاعات، واری داده ها و حتی لایروبی داده هاست که در حقیقت کشف دانش در پایگاه داده ها را توصیف می کند. بنابراین ایده ایی که مبنای داده کاوی است یک فرایند با اهمیت از شناخت الگوهای بالقوه مفید و قابل درک داده هاست. داده کاوی یا به تعبیر دیگر کشف دانش در پایگاه داده ها، استخراج بدیهی اطلاعات بالقوه مفید از روی داده هایی است که قبلا ناشناخته مانده اند. این مطلب برخی از روش های فنی مانند خوشه بندی، خلاصه سازی داده ها و فراگیری قاعده های رده بندی، یافتن ارتباط شبکه ها، تحلیل تغییرات و کشف بی قاعدگی را شامل می شود (سعیدی، ۱۳۸۴).

وب کاوی اشاره به کلیه فعالیت های داده کاوی و فنون وابسته دارد که برای کشف خودکار و استخراج دانش از اسناد و خدمات وب به کار می رود (ول و رویاکرز، ۲۰۰۴). اطلاعات بسیار زیاد و ناهمگونی در محیط وب وجود دارد که باعث می شود کسب دانش موجود در محتوای صفحات وب مشکل تر می شود، بنابراین در چنین محیطی به کارگیری ابزارها و فنون داده کاوی برای کشف اطلاعات و دانش مرتبط ضروری است (سالاری؛ ۱۳۸۳).

انواع روش های وب کاوی

وب کاوی ممکن است به سه شیوه انجام گیرد؛ کاوش محتوای وب، کاوش ساختار وب و کاوش کاربری وب می باشد. کاوش محتویات وب؛ فرایندی است که برای بدست آوردن مدل یا دانش ارزشمند و بالقوه از محتویات اسناد، اطلاعات توصیفی و یا نتایج جستجو در صفحات وب که هم زمان می تواند برای بدست آوردن اطلاعات مفید از ساختار وب و ارتباطات پیوندها به کار گرفته می شود.

کاوش ساختار وب؛ فرایندی است که ارتباطات پیوندهای بین مسیر صفحات، ساختار اسناد و ساختار مسیر در آدرس اینترنتی را جستجو می کند. در فضای وب علاوه بر محتویات صفحات، ساختار آن ها نیز اطلاعات مفیدی دارد، اگر ارجاعات به پیوندی افزایش یابد، پس آن صفحه مهم است و باید مسیر جستجو را بر اساس آن تغییر داد.

کاوش کاربری وب: فرایندی است که الگوی دسترسی به صفحات را با بررسی لاگ و داده های مربوط به آن کشف می کند. لاگ های سروری که همان اطلاعات تولید شده در هر ارتباط کاربر و سرور در جهان ثبت می شود. با بررسی این اطلاعات می توان رفتار کاربر را شناخت و ساختار صفحه وب را با هدف شخصی سازی بهبود بخشید. [ژائو و همکاران؛ ۲۰۱۴].

خوشه بندی

خوشه بندی یکی از بهترین روش هایی است که برای مدل سازی داده ها ارائه شده است و اطاعات را بر اساس شباهت به خوشه ها تقسیم می نماید. قابلیت آن در ورود به فضای داده و تشخیص ساختار آن ها، خوشه بندی را یکی از ایده آل ترین مکانیزم ها برای کار با دنیای عظیم داده ها کرده است. در خوشه بندی بدون اینکه هیچ مدل از پیش معینی داشته باشیم به دنبال یافتن مدل های مشترکی هستیم که د داده ها وجود دارد. ایده ی این روش د دهه ۱۹۳۵ ارائه شد و در حال حاضر با پیشرفت ها و جهش های عظیمی که در آن پدید آمده است کاربردهای مختلفی پیدا کرده است.

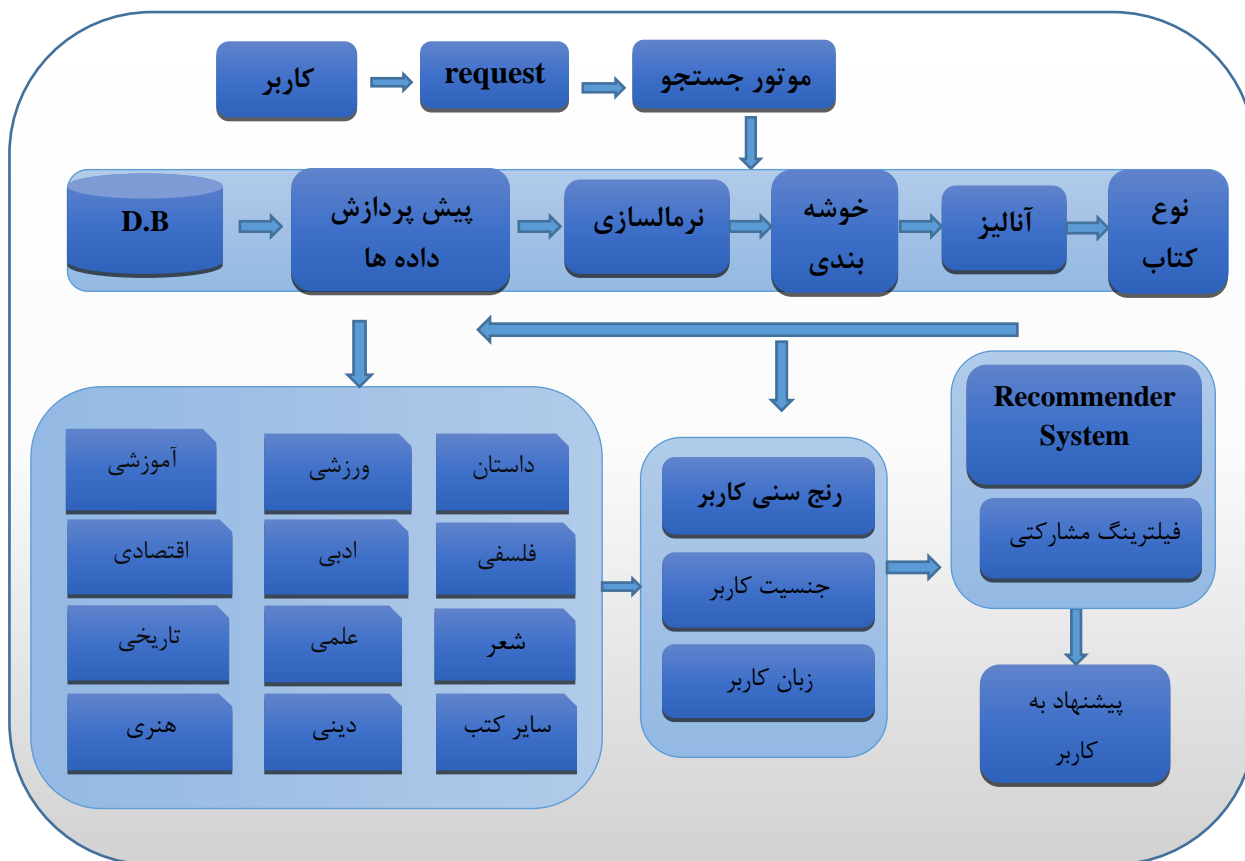
در سیستم های توصیه گر می توان از خوشه بندی برای دسته بندی کردن داده های مشابه با هم استفاده کرد تا جستجو و ساخت پیشنهادات فقط روی خوشه هایی انجام شود که پیش بینی می شود برای کاربر جالب تر خواهد بود. در نتیجه از تکنیک های خوشه بندی برای افزایش مقیاس پذیری الگوریتم های پشتیبان سیستم های توصیه گر استفاده شده است. این مساله به خصوص در شرایطی که برای پیشنهاد آیتمی لازم است تمام فضای مساله جستجو شود در کارایی الوریتم تاثیرات مثبتی نشان داده است [نکیتیم؛ ۲۰۰۹].

پرو فایل کاربر

در سال های اخیر تعداد زیادی سیستم توصیه گر برای پیش بینی علایق و ارائه توصیه به کاربران، طراحی و پیاده سازی شده است. در هر یک از این سیستم ها با توجه به حوزه کاری و اهداف، مجموعه ایی از تکنیک های ساخت، به روز رسانی و استخراج داده ها به کار گرفته شده است ولی محور اساسی در تمامی این سیستم ها پرو فایل کاربر است. چگونگی ساخت پرو فایلی که در ساخت توصیه ها استفاده خواهد شد، پرو فایل پیش فرض سیستم برای کاربران، نحوه بروز رسانی اطلاعات پرو فایل و منبع این بروز رسانی فاکتورهای هستند که در طراحی یک سیستم توصیه گر جایگاه مهمی دارند [مونتانا و همکاران؛ ۲۰۰۳].

روش پیشنهادی مورد نظر

در روش پیشنهادی مورد نظر پس از جمع آوری پایگاه داده ایی از کتب مختلف ابتدا عمل پیش پردازش داده ها را بر روی Data انجام می دهیم سپس باید Data را خوشه بندی کنیم برای این کار از روش خوشه بندی K-means استفاده می کنیم. پس از خوشه بندی Data باید داده ها را بر اساس نوع کتاب ها به دسته های مختلف (آموزشی، تاریخی، درسی، اقتصادی، ورزشی، سیاسی و...) دسته بندی می کنیم. سپس با اطلاعاتی که کاربر در پرو فایل خود وارد می نماید، سیستم توصیه گر مبتنی بر فیلترینگ مشارکتی بر اساس داده های موجود و آنالیزهای انجام شده بهترین تصمیم گیری را انجام می دهد و بهترین پیشنهاد را به کاربر ارائه می دهد که این پیشنهادات می توانند مورد علاقه کاربر باشند. در ادامه به شرح هر یک از مراحل این روش پرداخته می شود. در شکل (۱) نمایی از روش پیشنهادی را مشاهده می کنید.



شکل (۱) چارچوب پیشنهادی روش موردنظر

پیش پردازش داده‌ها

در مرحله اول روش پیشنهادی ابتدا باید عمل پیش پردازش داده‌ها را انجام دهیم در فرایند داده کاوی مانند طبقه بندی و خوشه بندی نیاز داریم تا داده‌ها برای الگوریتم آماده شوند، زیرا معمولا نمی‌توان داده‌ها را به صورت خام به الگوریتم‌های داده کاوی و یادگیری ماشین تزریق کرد. برای آماده سازی داده‌ها، نیاز است تا آنها را از شکل و حالت اولیه، خارج کرده و به شکلی که برای الگوریتم مناسب باشد تبدیل کرد. همچنین داده‌های موجود معمولا دارای زواید مختلفی هستند که ممکن است الگوریتم را دچار خطا کنند. در داده کاوی هم نیاز داریم تا داده‌های اضافی که به مسئله و الگوریتم کمکی نمی‌کنند را حذف کنیم. عملیات پیش پردازش داده‌ها معمولا قبل از عملیات اصلی الگوریتم‌های داده کاوی انجام می‌گیرند و باعث تسهیل و کمک به الگوریتم‌ها می‌شوند. پیش پردازش داده‌ها گامی مهم در راستای داده کاوی موفقیت آمیز است. اعمالی که در آماده سازی داده‌ها انجام می‌شود عبارتند از پاکسازی داده‌ها، یکپارچه سازی داده‌ها، تبدیل داده‌ها و کاهش داده‌ها می‌باشد. بر اساس نوع کاربردی که عمل داده کاوی باید روی آن انجام شود، تکنیک‌های مختلفی برای هر یک از این اعمال صورت می‌گیرد.

نرمال سازی داده‌ها

پس از پیش پردازش داده‌ها باید داده‌ها را نرمال کنیم، نرمال سازی داده‌ها تغییر داده‌ها به گونه‌ای است که آنها را به یک دامنه کوچک و معین مانند فاصله بین ۱- و ۱ نگاشت کنند. هدف نرمال سازی حذف افزونگی داد و باقی نگه داشتن وابستگی بین داده‌های مرتبط می‌باشد. این فرایند اغلب باعث ایجاد جداول بیشتر میشود ولی اندازه گیری پایگاه داده را کاهش داده و بهبود کارایی را تضمین می‌کند. روش‌های مختلفی جهت نرمال سازی داده‌ها وجود دارد که معروفترین آنها می‌توان به

روش MINMAXNORMALIZATIN اشاره کرد. در این روش هرکدام از داده ها را می توان به یک بازه دلخواه تبدیل کرد. فرمول کلی این روش برای تبدیل داده ها به بازه ی بین ۰ تا ۱ به صورت رابطه (۱) می باشد:

رابطه(۱)

$$Z = \frac{X - \text{MIN}(X)}{\text{MAX}(X) - \text{MIN}}$$

خوشه بندی داده ها

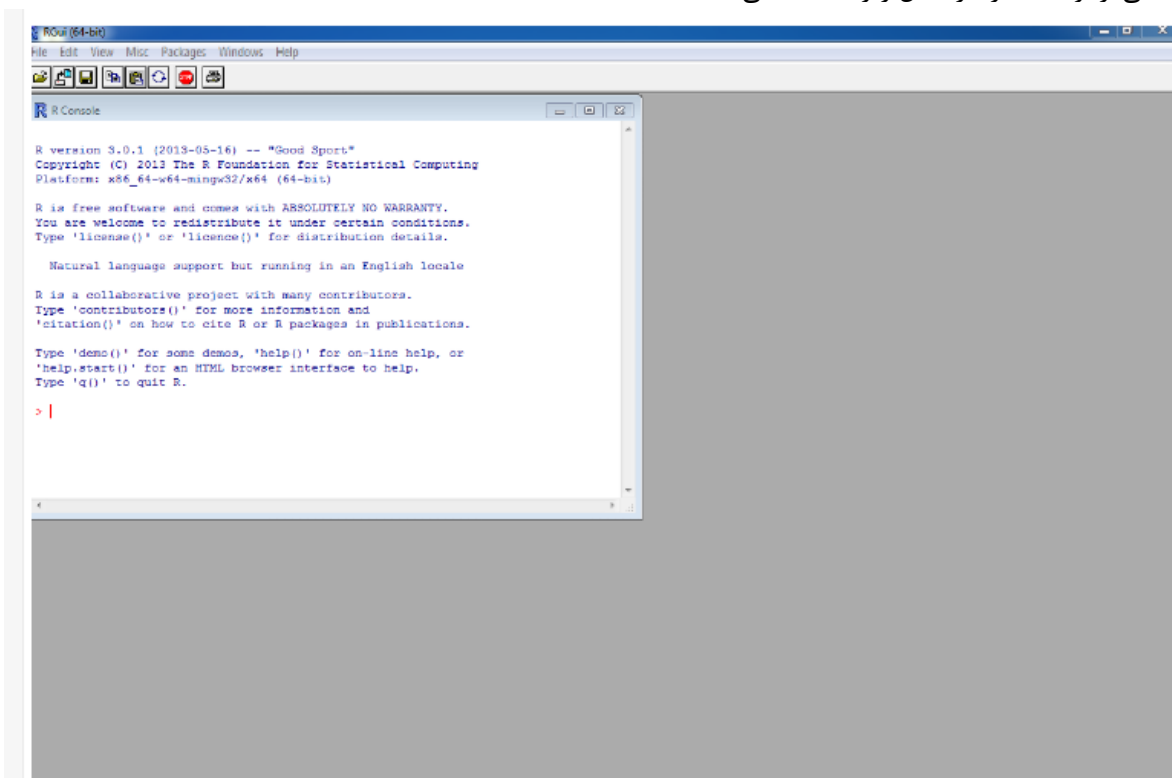
در این تحقیق از الگوریتم k-means به عنوان الگوریتم خوشه بندی استفاده می شود، این الگوریتم یکی از رایج ترین الگوریتم های استفاده شده در سیستم های توصیه گر می باشد. به عنوان معیار شباهت نیز ضریب همبستگی پیرسون استفاده می شود. در واقع در این بخش تعداد داده ها خوشه بندی می شود و برای هر خوشه یک نماینده و یا یک بردار میانگین در نظر گرفته می شود بدین صورت که تعداد کتب دیجیتالی به عنوان ورودی به الگوریتم خوشه بندی داده می شود و الگوریتم با استفاده ضریب همبستگی پیرسون کتاب ها را به خوشه هایی تقسیم می کند تا در هر خوشه بیشترین شباهت و در بین خوشه ها کمترین شباهت وجود داشته باشد. ضریب همبستگی پیرسون به صورت رابطه (۲) محاسبه می شود:

رابطه(۲)

$$SC(u, v) = \frac{\sum_j^m (r_{u,j} - \bar{r}_u)(r_{v,j} - \bar{r}_v)}{\sqrt{\sum_j^m (r_{u,j} - \bar{r}_u)^2 (r_{v,j} - \bar{r}_v)^2}}$$

آنالیز Data

پس از خوشه بندی داده ها نیاز داریم تا آنها را بر اساس انواع کتاب های مختلف آنالیز کنیم. نرم افزارهای زیادی برای آنالیز اطلاعات خوشه بندی شده موجود می باشد، در این تحقیق برای آنالیز داده های خوشه بندی شده از نرم افزار R استفاده می کنیم و توسط آن کتاب ها را به دسته بندی های مختلف دسته بندی می کنیم. نمایی از برنامه R را در شکل زیر مشاهده می کنید:



خصوصیات کاربر

در این مرحله با توجه به ویژگی‌هایی مانند رنج سنی کاربر، جنسیت کاربر، زبان کاربر و... سیستم توصیه گر سعی خواهد کرد که پیشنهادات خود را به کاربر ارائه دهد.

فیلترینگ مشارکتی

پالایش مشارکتی فرایند پالایش یا ارزیابی آیتم‌ها با استفاده از نظرات کاربران است. با وجود اینکه از شکل‌گیری اصطلاح پالایش مشارکتی تنها کمی بیش از یک دهه می‌گذرد، اما فلسفه این روش این است که استفاده از نظرات دیگران در تصمیم‌گیری است. برای قرن‌ها مورد استفاده‌ی انسان‌ها بوده است. برای مثال اگر دوستان شما از کتابی تعریف کنند راغب به خواندن آن کتاب خواهید شد. یا برعکس اگر کتابی را بد توصیف کنید بعید است برای خرید آن کتاب اقدام کنید. بعلاوه بعد از مدتی شما خواهید فهمید که نظرات کدام دوستانتان به نظرات شما نزدیکتر است و بتدریج فقط به نظرات آن دسته از دوستانتان که به شما شباهت دارند توجه خواهید کرد. کامپیوتر و بستر وب این امکان را فراهم کردند که قدم از ارتباطات شخصی فراتر گذاریم و به جای تصمیم‌گیری بر مبنای ده یا صد نفر از دوستان و نزدیکان، از نظرات بیش از میلیون‌ها کاربر استفاده کنیم. سرعت کامپیوتر این اجازه را می‌دهد که این نظرات را به صورت بلادرنگ پردازش کرده و بدانیم که افرادی که به ما شبیه هستند در مورد یک آیتم خاص که از آن بی‌اطلاعیم چه نظری دارند.

نحوه امتیاز دهی

در این قسمت به مفهوم امتیاز و ماتریس امتیاز پرداخته می‌شود. از آنجایی که در تکنیک پالایش مشارکتی نظرات کاربران نقش کلیدی دارد لازم است روش‌ها و قالب‌هایی برای جمع‌آوری آن طراحی شود. در ادبیات سیستم‌های توصیه گر روش‌های مختلفی برای جمع‌آوری نظرات کاربران معرفی شدند، اما روش معمولی که اکثر سیستم‌های توصیه گر مبتنی بر پالایش به کاربر گرفته می‌شود در نظر گرفتن یک بازه عددی (مثلاً ۱ تا ۵) برای هر آیتم، تعریف مفهوم هر یک از این اعداد (مثلاً ۱: بسیار بد، ۲: بد، ۳: متوسط، ۴: خوب، ۵: بسیار خوب) و درخواست از کاربر برای نگاشت یکی از این اعداد به هر یک از آیتم‌هایی که مشاهده می‌کند. این اعداد در ادبیات سیستم‌های توصیه گر امتیاز خوانده می‌شود و این روش امتیازدهی نام دارد. در جدول زیر ماتریس امتیاز یک سیستم توصیه گر مبتنی بر فیلترینگ مشارکتی چند کتاب را مشاهده می‌کنید:

امتیاز دهی بدین صورت می‌باشد که به عنوان مثال کاربر ۳ به کتاب بیگانه امتیاز ۳ و کاربر ۱ به کتاب قلعه حیوانات امتیاز نداده است.

شناسه کاربر	کتاب قلعه حیوانات	کتاب رهایی از زندان ذهن	کتاب ۱۹۸۴	کتاب بیگانه
کاربر ۱		۵	۲	۳
کاربر ۲	۴		۳	
کاربر ۳			۵	۳
کاربر ۴	۵	۵		۴

با استفاده از سیستم توصیه گر مبتنی بر پالایش مشارکتی کتاب‌هایی را به کاربر توصیه می‌کند که می‌تواند مورد علاقه کاربر باشد. برای پیش‌بینی کتاب‌ها، از تکنیک‌های CF انجام می‌شود، تکنیک‌های CF از یک دیتابیس از ترجیحات برای آیتم‌های توسط کاربران استفاده می‌کنند که برای تخمین موضوعات اضافی یا ایجاد کاربران جدید است. در یک سناریوی معمولی CF، لیستی از m کاربر وجود دارد $\{u_1, u_2, u_3, \dots, u_m\}$ و همچنین یک لیست از n آیتم $\{i_1, i_2, i_3, \dots, i_n\}$ و هر کاربر U_i یک لیست از آیتم‌ها دارد l_{ui} که کاربر آنها را رتبه‌بندی کرده است، و یا آنهایی که ترجیحات آنها از طریق رفتارشان استنباط کرده است. رتبه می‌تواند اشاره‌های صریح باشد و... که روی مقیاس ۱ تا ۵ است و یا همچنین می‌تواند اشاره ضمنی باشد.

الگوریتم نزدیکترین همسایه مبتنی بر کاربر

این نوع از الگوریتم‌ها در پیش بینی علاقه ی یک کاربر به یک آیتم خاص از امتیازهایی که کاربران مشابه او به آن آیتم داده اند استفاده می کند. این کاربران مشابه، اصطلاحاً همسایگان کاربر نامیده می شوند. اگر n کاربر مشابه کاربر u باشد، گفته می شود که n همسایه u است. برای پیش بینی میزان علاقه کاربر u به آیتم i باید میانگین امتیازهایی که همسایگان u (شامل کاربر n) به i داده اند را با تساوی (۱،۲) حساب کرد. در این تساوی r_{ni} امتیازی است که کاربر n به آیتم i داده است. نحوه محاسبه این الگوریتم را در رابطه (۲) مشاهده می کنیم.

رابطه (۳)

$$pred(u, i) = \frac{\sum_{n \in neighbors(u)} r_{ni}}{\text{number of neighbors}}$$

پیشنهاد به کاربر

در مرحله آخر سیستم توصیه گر با تمام آنالیزهایی که انجام داده، لیستی از کتاب های مختلف را به کاربر پیشنهاد خواهد کرد که این کتاب ها می تواند مورد علاقه کاربر باشد.

نتیجه گیری

پیشرفت ها و تحولات پیش آمده در کتابخانه های دیجیتالی به عنوان نسل جدیدی از کتابخانه ها، ضرورت ارزیابی قابلیت های جستجو و نمایش نتایج جستجو را ایجاد می نماید، چرا که نظام جستجو و بازیابی هر کتابخانه دیجیتالی باید پاسخگوی تمام گروه های کاربری خود باشد، یعنی با استفاده از آن هم جستجوگران مبتدی و کم تجربه و هم جستجوگران با تجربه بتوانند به مطلوب خود برسند و با فراهم ساختن امکانات جستجوی متعدد اجازه انجام انواع مختلف جستجوها را به کاربران دارای سطوح مهارتی مختلف بدهند. در این تحقیق سعی شد با توجه به اینکه کتب مختلف روزانه در حال افزایش می باشند و دسترسی به این کتب با مشکلاتی از قبیل نحوه دسترسی، هزینه و... مواجه شده است به طراحی یک سیستم توصیه گر مبتنی بر پالایش مشارکتی برای استخراج کتب موردنیاز کاربران از کتابخانه دیجیتالی پرداخته شد که می تواند کتب موردنظر کاربر را از حجم وسیع اطلاعات موجود در کتابخانه دیجیتالی استخراج و بتواند پیشنهادات مناسب را به کاربر ارائه دهد.

منابع و مراجع

- [۱] پناهی، سمیه (۱۳۹۶). پژوهشی در باب خدمات شخصی سازی کتابخانه دیجیتال براساس داده کاوی. مجله الکترونیکی. متبادار ۰/۲ دوره سوم. شماره اول (اردیبهشت ماه).
- [۲] رحمانی، محمود و حاجی زین العابدینی، محسن، (۱۳۹۵). کاربردهای داده کاوی در علم اطلاعات و دانش شناسی: نشریه مدیریت اطلاعات و دانش شناسی. مقاله ۲، دوره ۲، شماره ۳، پاییز ۱۳۹۴، صفحه ۲۳-۳۲.
- [۳] سالاری، عبدالرضا. (۱۳۸۳) ”کاربرد محاسبات نرم در وب کاوی”. پایان نامه کارشناسی ارشد کامپیوتر، دانشگاه آزاد اسلامی واحد علوم و تحقیقات.
- [۴] سعیدی، احمد. (۱۳۸۴) ”داده کاوی، مفهوم و کاربرد آن در آموزش عالی”. شماره ۱۸، اسفند، ص ۲-۱۷.
- [۵] کوشا، کیوان. ۱۳۸۴. کتاب رقومی چیست؟ اصطلاحی رایج با مفهوم ابهام برانگیز. مطالعات کتابداری و سازماندهی، ۹۷: ۶۳-۱۱۰.
- [۶] مرادی، گلمراد؛ قاسمی، وحید (۱۳۹۱) تکنیک داده کاوی و کاربرد آن در مطالعات اجتماعی. مجله علوم اجتماعی. بهار و تابستان، شماره ۱۱۹. (۱۷۶-۱۵۵).
- [۷] نوروزی، یعقوب. ۱۳۸۹. تحلیلی بر کاربر مداری رابط کاربر در صفحات وب فارسی کتابخانه های دیجیتالی ایران و ارائه پیشنهادی- فصلنامه علوم و فناوری اطلاعات، ۱۶(۸): ۷۹-۶۴.
- [8] Babu, P.B. Krishnamurthy, M. 2013. Library Automation to Resource Discovery: A Review of Emerging Challenges. *Electronic Library*, 31(4):433-451.
- [9] Burke, R., —Hybrid recommender systems: Survey and Experiments. || *User Model. UserAdapt. Interact.* 12, 4, 331–370, (2002)
- [10] Kim, H. Lee, J. Y. 2008. Exploring the emerging intellectual structure of archival studies using text mining: 2001-2004. *Journal of Information Science*, 34(3): 356-369.
- [11] Konstan, J. A., —Introduction to recommender systems. || In *Proceedings of the 2008 ACM SIGMOD international Conference on Management of Data*, Vancouver, Canada, (Jun, 2008).
- [12] Lu Jie, Dianshuang Wu, Mingsong Mao, Wei Wang, Guangquan Zhang, (2015). *Recommender System Application Developments: A Survey*, *Decision Support Systems*, Volume 74, 12-32.
- [13] Montaner, M., Lopez, B., Rosa, and J.L.D.L., taxonomy of recommender agents on the internet. || *Artificial Intelligence Review* 19 (2003).
- [14] Nkweteyim, D., —Hyperlink Recommender Systems Design: a Research Study on Tools and Techniques. || *VDM Verlag*. (2009)
- [15] Pareek, S. 2013. A Bibliometric analysis of the literature of IFLA Journal during 2001-2010. <http://digitalcommons.unl.edu/libphilprac/954>. (accessed 12/07/ 2013)
- [16] Queiroz da Silva Edjalma, Celso G. Camilo-Junior, Luiz Mario L. Pascoal, Thierson C. Rosa, (2016). An evolutionary approach for combining results of recommender systems techniques based on collaborative filtering, *Expert Systems with Applications*, Volume 53, 204-218.
- [17] R. Burke, —Knowledge-Based Recommender Systems, || *Encyclopedia of Library and Information Systems*, A. Kent, ed., vol. 69, Supplement 32, Marcel Dekker, (2000).
- [18] Rao Kagita Venkateswara, Arun K. Pujari, Vineet Padmanabhan, (2015). Virtual user approach for group recommender systems using precedence relations, *Information Sciences*, Volume 294, 15-30.
- [19] Saracevic, T. 2000. Digital library evaluation: Toward evolution of concepts. *Library trends*, 49 (2):350-369. Shen, R. Gonçalves, M.A. Fox, E.A.. 2013. Key Issues Regarding Digital Libraries: Evaluation and Integration. *Synthesis Lectures on Information Concepts, Retrieval, and Services*, 5 (2):1-110. Tedd
- [20] Wang, H., Li, W.-J., (2014). Relational collaborative topic regression for recommender systems, *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, Volume 27, Issue 5, 1343 - 1355.

- [21] Wang, H., Li, W.-J.,(2014). Relational collaborative topic regression for recommender systems, IEEE Transactions on Knowledge and Data Engineering (TKDE), Volume 27, Issue 5,1343 - 1355.
- [22] Wei Jian, Jianhua He, Kai Chen, Yi Zhou, Zuoyin Tang, (2016). Collaborative Filtering and Deep Learning Based Recommendation System for Cold Start Items, Expert Systems with Applications, Volume 69, 29-39.
- [23] Wel, litavan; Royackers, Lamber(2004)"Ethical issue in web data mining". Ethics and Information Thecnology,6, pp. 129 140
- [24] Yang Xiwang, Yang Guo, Yong Liu, Harald Steck,(2014). A survey of collaborative filtering based social recommender systems, Computer Communications, Volume 41,1-10.
- [25] Zhao, L. Zhang, Q.2011. Mapping knowledge domains of Chinese digital library research output, 1994–2010. Journal of scientometrics, 89: 51– 87.
- [26] Zhao, Yonghua, and Hong Lin. "WEB data mining applications in e-commerce.", IEEE 9th International Conference on Computer Science & Education, Vancouver, Canada, 2014.