

Web-based Personalized Hybrid Book Recommendation System

۱ مقدمه

در این مقاله سعی شده مدلی از ركامندر سيستم هايبريد بر پايه وب براي معرفي كتاب، نشان داده شود. اين سيستم با بررسي ابعاد مختلف پيشنهاد كتاب، كارايي بهتر از دو سيستم تك بعدي content based و collaborative دارد.

۲ سيستم پيشنهاد شده

اين سيستم از ۳ بخش سرور مركزي، و انبار داده و اينترنت تشكيل شده است.

در سرور مركزي تمام عمليات هاي محاسباتي انجام خواهند شد. اين سرور به يك انبار داده اي متصل است كه شامل ۲ جدول است. در جدول اول اطلاعات يوزر و ويژگي هاي demographic مربوط به آن، در جدول دوم اطلاعات كتاب ها مانند اسم، عكس و ...، در جدول سوم اطلاعات بدست آمده از web scraping و در جدول چهارم امتياز هاي افراد به كتاب ها قرار دارند. جدول چهارم شامل ويژگي مهمي به نام timestamp هست كه به سيستم كمك ميكند تا كتاب هاي قديمي را از دور خارج كند و همگام با تغيير سلايق مردم پيش برود.

در بک اند این سیستم میتوان از Apache Mahout و يا MySQL استفاده کرد. در صورت استفاده از Mahout، با استفاده از ويژگي Hadoop آن، ميتوان در ابعاد بزرگ تري نسبت به MySQL كار كرد اما سرعت دسترسي آن نسبت به MySQL بسيار كمتر است.

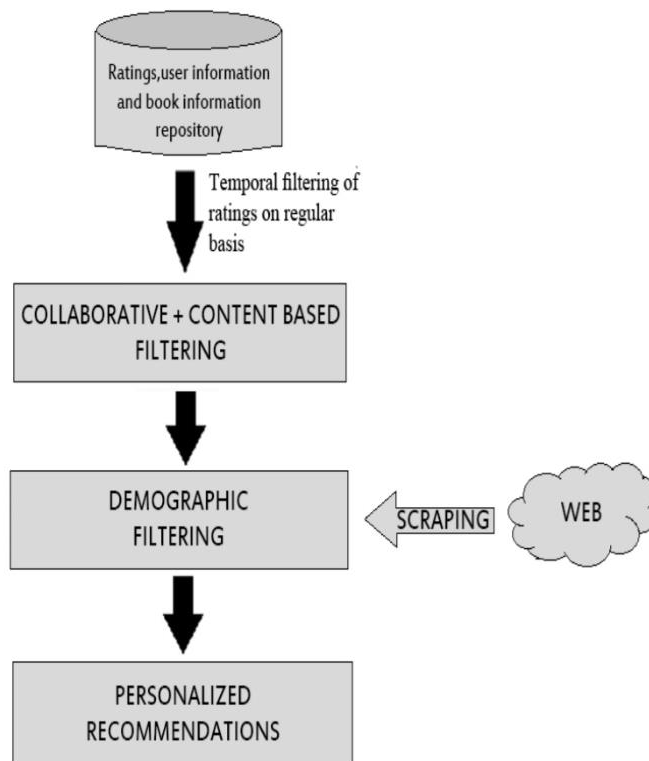
از متدهاي فيلترينگ demographic، collaborative item-item، user-user و content-based در ساخت اين سيستم استفاده شده.

ميتوانيد در شكل ۱-۲، ساختار و روند كار اين سيستم را مشاهده نماييد.

۳ متدهاي فيلترينگ

در اين سيستم، از تركيبی از ۳ متد فیلترینگ collaborative، content based و demographic استفاده شده.

ابتدا به معرفي اين ۳ متد در سيستم پيشنهاد شده ميپردازيم.



شکل ۱-۲: روند کار سیستم پیشنهادی

۱-۳ Collaborative-based filtering

در Apache Mahout کلاس‌های `PearsonCorrelationSimilarity` و `GenericUserBasedRecommender` وجود دارند که در پیاده‌سازی مقادیر شباهت `user-user` و `item-item` به ما کمک میکنند.

۲-۳ Content-based filtering

در این متد علاوه بر نمرات داده‌شده توسط کاربران، از اطلاعات دیگری مانند نویسنده، ژانر، هزینه، تعداد صفحات و ... نیز برای پیدا کردن بهترین پیشنهادات استفاده میشود.

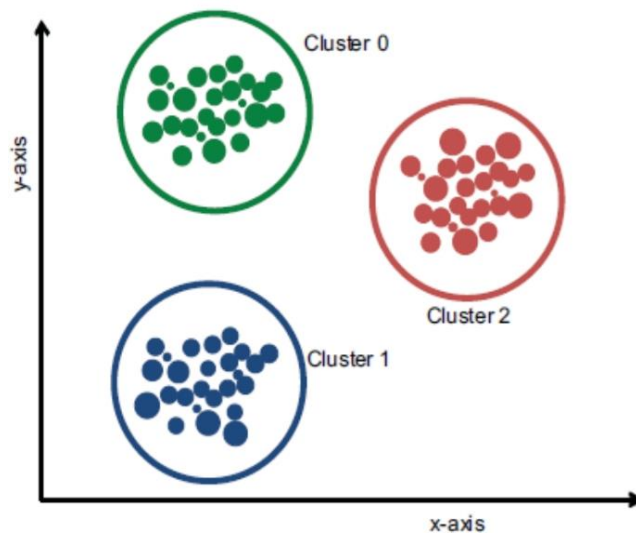
مهمترین قدم در این بخش، پیدا کردن متد یادگیری است زیرا علاوه بر بازده، به دلیل بزرگی داده‌ها، باید به فاکتورهایی مانند زمان و فضا توجه کرد.

در Apache Mahout، کلاسی برای پیاده‌سازی `Content-based filtering` وجود ندارد اما کلاس‌هایی وجود دارند که میتوانند برای پیاده‌سازی الگوریتم‌های موردنیاز، به ما کمک کنند.

۳-۳ Demographic-based filtering

تکنیک کلاسترینگ میتواند در جهت افزایش بهره‌وری، به سیستم کمک کند. در این تکنیک، به جای استفاده از همه یوزرها به عنوان training set، یوزرها را بر اساس ویژگی‌هایی مانند سن، جنسیت و ... به دسته‌هایی تبدیل میکند و سپس از این دسته‌ها به عنوان training set استفاده میشود. در شکل ۳-۱ میتوانید ساختار گرافیکی این تکنیک را ببینید.

در Apache Mahout، کلاس‌هایی مانند SequenceFile، EuclideanDistanceMeasure و ManhattanDistanceMeasure وجود دارند که به ما در پیاده‌سازی این متد کمک میکنند.

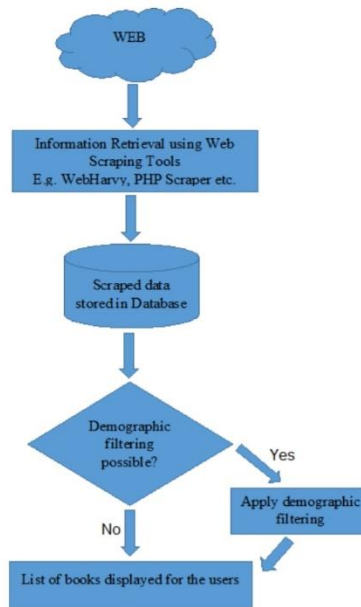


شکل ۳-۱: تکنیک کلاسترینگ

۴ Web Scraping

به عمل جمع‌آوری اطلاعات از سطح اینترنت Web Scraping گفته میشود. در شکل ۴-۱ میتوانید فلوجارت عملیات Web Scraping را مشاهده نمایید.

وبسایت‌های فعال در حوزه E-Commerce میتوانند در سیستم، به ما کمک زیادی بکنند. به عنوان مثال با انجام Web Scraping بر روی وبسایت آمازون، میتوان اطلاعاتی از قبیل پرفروش‌ترین کتاب‌ها، کتاب‌ها با بیشترین امتیاز و ... را به دست آورد و برای حل مشکلاتی از قبیل cold start و gray sheep users، از آنها استفاده کرد.



شکل ۴-۱: فلوچارت web scraping

۴ اعتماد

این بخش یکی از مهمترین بخش‌های سیستم پیشنهادی است زیرا در صورتی که کاربر به سیستم اعتماد نداشته باشد، دیگر ابعاد سیستم بی‌فایده خواهند بود. برای این موضوع میتوان در صفحه‌ای از وبسایتان، متن یا ویدیویی با موضوع روش کار این سیستم در اختیار کاربران بگذاریم تا با درک روند کار سیستم، اعتمادشان افزایش یابد.

روش دیگر استفاده از mirror behavior است. در این روش، کاربران میتوانند پروفایل‌های دیگر کاربران را مشاهده نمایند و کاربران مشابه خود را تشخیص دهند. این روش باید بصورت ۲ طرفه انجام شود و در صورتی که کاربری اجازه مشاهده پروفایلش را به دیگران ندهد، خود نیز دسترسی به پروفایل دیگران نداشته باشد.

در صورت افزایش امنیت سیستم، کاربران به سیستم راحتتر اعتماد میکنند. برای جلوگیری از حملاتی مانند نظرها و امتیازات جعلی، میتوان از متدهایی مانند کپچا، بررسی آیدی کاربر و ... استفاده کرد.