

پارتیشن‌بندی زمانی گراف برای خوشه‌بندی در سیستم‌های برجسب دار

علی‌اکبر اله‌داغی^{۱*}، مهرداد جلالی^۲، سید جواد سید مهدوی‌چابک^۳

^{۱،۲،۳} دانشگاه آزاد اسلامی، واحد مشهد، گروه کامپیوتر، مشهد، ایران

چکیده:

در عصر حاضر رشد اطلاعات در دنیای وب ۲.۰ به حدی بوده است که به علت حجم وسیع داده‌ها و اطلاعات و همچنین تغییر برخی مفاهیم در طی گذر زمان، اطلاعات غیرضروری و غیر مرتبط بسیاری با آنچه کاربران به دنبال آن هستند، به وجود آمده است. در این مقاله ما برای حل این مشکل، خوشه‌بندی زمانی برجسب‌ها در سیستم‌هایی که از برجسب به‌عنوان یک متاداده استفاده کرده‌اند در طی زمان در حال تغییر می‌باشند را پیشنهاد می‌دهیم. این عمل از طریق پارتیشن‌بندی زمانی گراف برجسب‌هاست، به طوری که با تغییر دادن وزن‌های مشابهت برجسب‌ها در طی زمان، خوشه‌بندی ما هم دچار تغییر شده و خودش را با تغییرات وفق می‌دهد. برای نشان دادن کارایی این روش، ما آن را بر روی مجموعه داده سایت متافیلتر اجرا کرده و با روش‌های مشابه مقایسه کردیم. نتایج نشان‌دهنده این مطلب است که روش پیشنهادی ما، F-Measure میانگین خوشه‌ها را ۲۴٪ نسبت به بهترین روش مشابه در طول زمان، بهبود بخشیده است و از نظر مفهومی، هم با مفاهیم گذشته هم با مطالب روز در ارتباط است.

واژه‌های کلیدی:

برجسب، خوشه‌بندی، Temporal, Graph Partitioning, Folksonomy

۱- مقدمه

در عصر حاضر تبادل و رشد دانش به کمک فناوری اطلاعات در حال افزایش نمایی است. یکی از مهم ترین عوامل این افزایش اشتراک اطلاعات بین کاربران وب است. این اشتراک گذاری در زمان های اخیر تسهیل و افزایش یافته است که در اصطلاح، وب ۲.۰ نامیده می شود [۱]. وب ۲.۰ بر پایه چند اصل اساسی بنا شده است. از این میان سه اصل زیر ماهیت وب ۲.۰ را که برخلاف فناوری های دیجیتال و وب ۱.۰ است را تشکیل می دهد. ۱. وب یک جایگاه^۱ است. ۲. کاربر داده های خود را کنترل می کند. ۳. کاربر نقش مشارکتی دارد [۲]. ترکیب این سه اصل تولید برنامه های مبتنی بر وب را به شیوه ای ممکن ساخته است که وب را به یک وب مشارکتی تبدیل می کند.

وب ۲.۰ جنبه های متعددی دارد که نمونه های معمولی از مشارکت کاربران را انعکاس می دهد. این جنبه ها عبارت اند از: ویلاگ نویسی^۲، برچسب زدن^۳، نشانه گذاری های اجتماعی^۴، شبکه های اجتماعی^۵، پادکست و غیره [۱]. در این میان، برچسب ها به عنوان یک متاداده^۷ در نمایش خلاصه محتوای هر سند بسیار حائز اهمیت هستند.

برچسب ها در ضمن سادگی، توسط کاربران برای توصیف و یا حاشیه نویسی از هر نوع منبع و یا سندی برای بازیابی در آینده اختصاص داده شده اند. برچسب های ارسال شده توسط گروهی از کاربران Folksonomy را تشکیل می دهد [۳]. هدف از برچسب گذاری در یک فکسونومی این است که جستجو، حرکت بین اطلاعات و درک یک بدنه اطلاعات، در طول زمان و به تدریج برای

کاربران آسان شود. از آنجا که فکسونومی ها عموماً در محیط های مبتنی بر اینترنت به وجود آمده اند، کاربران می توانند دریابند که یک برچسب معین را چه کسی استفاده می کند و همچنین او از چه برچسب های دیگری استفاده می کند. به این ترتیب کاربران می توانند به درک متقابلی از مجموعه برچسب هایی برسند که دیگران برای تفسیر، درونی سازی و درک یک بخش از اطلاعات نیاز دارند. این مورد می تواند در سیستم های مدیریت دانشانسان محور یک گام بزرگ روبه جلو باشد؛ چرا که قابلیت افراد در ساختاردهی صحیح دانش و همچنین یافتن محتواهای مرتبط را بالا می برد. این اصل، دریافتن اطلاعات مرتبط^۸ نام دارد. بخش عمده ای از جاذبه فکسونومی در توانمندی ذاتی و راحتی آن دریافتن اطلاعات نسبت به ابزارهای سنتی جستجو مانند موتورهای جستجو است [۴].

ساختار این مقاله بدین صورت است:

بخش ۲. مرور کارهای انجام شده در زمینه

خوشه بندی متون و سیستم های برچسب دار: در این بخش کارهای گذشتگان را که مرتبط با زمینه خوشه بندی اطلاعات صفحات وب با رویکرد استفاده از برچسب ها است را بررسی کرده و مورد تجزیه و تحلیل قرار می دهیم.

بخش ۳. روش پیشنهادی:

در این بخش روش پیشنهادی خود را مطرح کرده و شیوه پیاده سازی و الگوریتم های موجود و نوآوری خود را شرح خواهیم داد.

بخش ۴. نتایج و ارزیابی:

در این بخش نتایج و مقایسه های خود را با روش های موجود مشابه بررسی و ارزیابی خواهیم کرد.

بخش ۵: نتیجه گیری و پیشنهادها:

نتیجه گیری ها و پیشنهاد های خود اشاره خواهیم نمود.

1- Platform
2- Blogging
3- Tagging
4- Social Bookmarking
5- Social Networking
6- Tag
7- Metadata

8- Pivot Browsing

۲- مروری بر کارهای گذشته

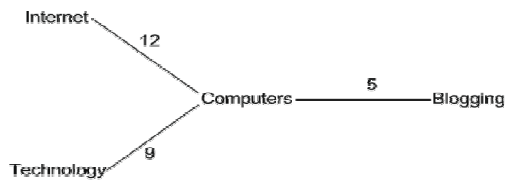
از جمله مقالاتی که از برچسب به عنوان داده برای خوشه‌بندی خود استفاده کرده‌اند، می‌توان به کارهای [۳، ۵، ۶، ۷، ۸، ۹، ۱۰، ۱۱، ۱۲] اشاره نمود که آنها را در ۳ قالب ذیل تقسیم‌بندی می‌کنیم:

۱. کارهایی که مربوط به خوشه‌بندی برچسب‌ها و یا استفاده از برچسب‌ها برای دسته‌بندی کردن اسناد و یا کمک به موتورهای جستجو می‌باشند، مانند [۳، ۵، ۷ و ۸].
۲. مجموعه کارهایی که علاوه بر خوشه‌بندی برچسب‌ها به پارامتر زمان آن‌هم توجه داشته‌اند و در خوشه‌بندی از آن متأثر بوده‌اند، از جمله مقالات و کارهای [۳، ۵ و ۷].
۳. کارهایی که مربوط به خوشه‌بندی اسناد متنی است؛ مانند [۵، ۱۱ و ۱۲] که در این بین مقالات [۵ و ۱۱] از برچسب‌ها نیز کمک گرفته و یا به نوعی کار خود را با آن مقایسه کرده‌اند.

۱-۲- خوشه‌بندی برچسب‌ها

در [۵] نویسندگان ابتدا به تفاوت‌های بین وب ۱.۰ و وب ۲.۰ اشاره کردند و خاطر نشان کردند که در وب ۲.۰ خلاف وب ۱.۰- که اکثراً خدمات دهنده‌ها مطالب را در اختیار کاربران می‌گذاشتند- این کاربران هستند که خود برای سایر کاربران مطلب تولید می‌کنند. در ادامه کارهای مرتبط با خودشان را مقایسه کردند و اشاره به این موضوع داشتند که اکثراً در حجم وسیع دچار نفرین ابعاد می‌شوند. کار اصلی این مقاله مبتنی بر برچسب‌های پست داخل وبلاگ- هاست؛ به طوری که بر اساس ارتباط بین برچسب‌های گراف وزنداری را ایجاد می‌کند و وزن یال‌ها در واقع برابر با تعداد پستی است که دو برچسب موجود در دو طرف یال را در بین برچسب‌های خود آورده است شکل (۱). سپس بر اساس پاریشن‌بندی گراف خوشه‌بندی خود را انجام می-

دهد و نام این کار را خرد جمعی می‌نامد. البته آنها این کار را با چند روش دیگر مورد بررسی قرار می‌دهند.



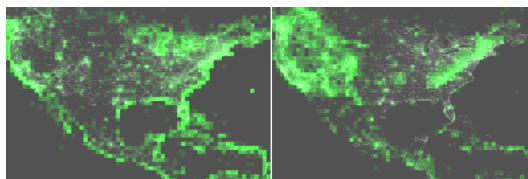
شکل (۱): یک نمونه از گراف ارتباط برچسب‌ها [۱]

در کار بعدی سوین انفینسن و همکاران در [۳]، به کاربرد ابزارهای وب ۲.۰ در کتابخانه‌های دیجیتال اشاره نمودند که باعث می‌شود مفهوم کتابخانه ۲.۰ پدید آید. آنها نشان دادند که با اشتراک‌گذاری عناوین کتاب‌ها و تعیین برچسب‌ها برای هر یک توسط اساتید، ناشرین و دانش‌آموزان چطور می‌توان به دسترسی سریع‌تر کاربران برای جستجوی یک کتاب خاص کمک کرد. در [۱۱] Brooks و Montanez تأثیر استفاده از برچسب‌گذاری توسط انسان را برای خوشه‌بندی پست‌های وبلاگ تحلیل کردند. آنها دریافتند که هر چند برچسب‌ها برای طبقه‌بندی پست‌های وبلاگ به دسته گسترده مؤثر هستند، ولی آنها کمتر در محتوای خاصی از پست، مؤثر نشان داده می‌شوند. سپس آنها استخراج یک مجموعه متاهی از بیشترین کلمات مهم از مقاله خودشان را به وسیله معیار *tf-idf* پیشنهاد دادند و ذکر کردند که استفاده از آن‌ها برای طبقه‌بندی بازدهی بهتری دارد.

۲-۲- خوشه‌بندی زمانی برچسب‌ها

هایپن ژانگ و همکارانش در [۶] از روی ارتباط برچسب‌ها با مکان و زمان، خوشه‌بندی‌های خود را انجام دادند. مجموعه‌هاده آنها شامل تصاویری است که هر تصویر

علاوه بر برچسب‌هایی برای معرفی آن تصویر، دارای عنوان، زمان ثبت و موقعیت جغرافیایی می‌باشد.



شکل (۲): توزیع جغرافیایی برای برچسب "beach" (تصویر چپ) و "mountains" (تصویر راست) [6]

آنها به ۵ روش مختلف برچسب‌های موجود را خوشه‌بندی می‌کنند.

۱. اولین روش آنها استفاده از توزیع جغرافیایی برچسب‌های تصاویر است. در شکل (۲) نمونه‌ای از توزیع جغرافیایی تصاویر گرفته‌شده با دو برچسب مختلف نشان داده شده است. آنها از روی این توزیع یک بردار ویژگی تحت عنوان بردار ویژگی Geo-spatial ایجاد می‌کنند.

۲. دومین روش ایجاد یک بردار ویژگی زمانی از روی توزیع بازه‌های زمانی استفاده از برچسب‌هاست.

۳. روش دیگر خوشه‌بندی آنها استفاده از دو خاصیت توزیع مکانی و زمانی، یعنی ترکیب هر دو روش قبلی است.

۴. استفاده از ویژگی وقوع مشترک زوج برچسب‌ها است. شیوه محاسبه آن هم استفاده از رابطه (۱) است:

$$\text{mutual_info}(t_1, t_2) = \log \left(\frac{\text{co_occur}(t_1, t_2)}{\text{occur}(t_1) \times \text{occur}(t_2)} \right) \quad (1)$$

۵. فقط استفاده از $\text{co_occur}(t_1, t_2)$

آنها برای ارزیابی کار خود سایتی طراحی کرده و از مردم دعوت کردند که به نتایج خوشه‌بندی‌های ایجادشده هر روش امتیاز بدهند؛ به عبارت دیگر، ارتباط برچسب‌ها را از

طریق نظرسنجی عمومی حاصل کردند. برای مقایسه هم از معیارهای متریک precision و recall استفاده کردند.

۲-۳- خوشه‌بندی‌های متنی

در [۱۲] نویسنده به بررسی چند روش مفید برای خوشه‌بندی وبلاگ‌ها با توجه به متن موجود در آنها می‌پردازد. این روش‌ها عبارتند از: LSI^1 , VSM^2 , K-Means و FCM^3 . کار اصلی این پروژه ارائه یک اصلاح الگوریتم بهینه انتخابی برای خوشه‌بندی وبلاگ است. در انتها در یک جدول الگوریتم‌های فوق را با همدیگر مقایسه می‌کند. در جمع‌آوری داده از تعدادی وبلاگ‌هایی که قبلاً به صورت آفلاین گرفته شده است، استفاده می‌کند. سپس در مرحله تبدیل‌ها ضمن حذف برچسب‌های HTML، کلمات ربطی مانند "of", "is", "on", "the" را از داده‌ها حذف می‌کند. همچنین معادل‌سازی را برای برخی کلمات که در جمله از شکل‌های دیگر آن استفاده شده است، انجام می‌دهد. مثلاً کلمات iron معادل irons است و یا به جای damaged از شکل ساده آن یعنی damage استفاده می‌کند.

برای وزن دهی به کلماتی که در متن ظاهر شده‌اند از روش TF-IDF استفاده می‌کند.

۳- روش پیشنهادی

در این بخش روش کار خود که پارتیشن‌بندی گراف زمانی برای خوشه‌بندی در سیستم‌های برچسب دار است را شرح می‌دهیم. قبل از شروع، ذکر این نکته لازم است که خوشه‌بندی ما در این سیستم‌ها مستقل از کاربر است و فقط با منابع و برچسب‌گذاری سایر کاربران ارتباط دارد و به

¹- Latent Semantic Indexing

²- Vector Space Model

³- Fuzzy C-Means

۳-۳- ماژول خوشه بندی برچسب ها

در خوشه بندی برچسب ها مطابق [5، 8 و 9] ابتدا روش مشابهت خود را تعیین می کنیم. ما روشی را انتخاب می کنیم که قرار است مشابهت بین دو برچسب را نشان دهد و در ضمن نرمالیزه نیز باشد؛ لذا از معادله (۲) استفاده می کنیم.

$$Sim(Tag_i, Tag_j) = \frac{|A \cap B|}{|A \cup B|} \quad (2)$$

در معادله (۲) مجموعه سندهایی است که دارای Tag_i و B مجموعه سندهایی است که دارای Tag_j می باشند. در ادامه برای انجام خوشه بندی، برچسب هایی را که قرار است در خوشه بندی شرکت کنند، انتخاب می کنیم. روش انتخاب هم بدین صورت است که بر اساس فراوانی استفاده از برچسب ها، آنها را به صورت نزولی مرتب کرده و K تاییابتدای آن را انتخاب می کنیم. پس از انتخاب برچسب ها جهت خوشه بندی، ماتریس مشابهت پایین مثلثی ای بر اساس معادله (۲) تشکیل می دهیم. روش خوشه بندی ما استفاده از پارتیشن بندی گراف است، به طوری که پس از تشکیل گراف از روی ماتریس فوق، عدد τ را به عنوان حد آستانه در نظر گرفته و یال های کمتر از آن حد را حذف می کنیم. حال به دو روش Single Link یا Clique می توانیم عمل پارتیشن بندی را روی گراف انجام دهیم.

۳-۴- ماژول سازمان دهی مجدد خوشه ها

تا این مرحله ما یک مرحله خوشه بندی کامل را انجام دادیم ولی پس از این نیاز است که خوشه ها به روز نگه داشته شوند که به روش ذیل عمل می کنیم:

۱. در هر بازه زمانی T کسر مشخص از یال های گراف مطابق معادله (۳) کم می کنیم. این عمل به علت آن است که

نوعی باخرد جمعی همراه است. ماژول های کاری به ترتیب ذیل خواهند بود:

۱. ماژول خزنده^۱ صفحه های وب (جمع آوری مستندات برچسب دار)

۲. ماژول آنالیز سند شامل: استخراج متن، پیش پردازش اولیه، استخراج برچسب ها

۳. ماژول خوشه بندی برچسب ها

۴. ماژول سازمان دهی مجدد خوشه ها

۳-۱- ماژول خزنده صفحه های وب

در ابتدا لیستی از سایت های مرجع تهیه می شود که در آن کاربران ضمن گذاشتن منبع (لینک، متن، تصویر و ...) برای آنها برچسب هایی را نیز قرار می دهند. سپس مشابه [5] خزنده ای طراحی شده است که با رجوع به ریشه تک تک سایت های مرجع، به صورت اول سطح شروع به خزیدن می کند.

۳-۲- ماژول آنالیز سند

ورودی این ماژول فایل حاوی تگ های HTML است و خروجی آن برچسب های موجود در آن صفحه. در این مرحله ما از پیش مکان قرار گرفتن برچسب ها را می دانیم و علت این است که قبلاً یک فرد خبره ساختار صفحات مرجع را بررسی کرده و از مکان های قرار گرفتن مشخصات منبع مطلع است. به همین جهت این ماژول با دریافت صفحه وب مورد نظر به سراغ تگ های html از پیش تعیین شده می رود. سپس نام و مشخصات آن منبع را به همراه برچسب های آن، در بانک اطلاعاتی با زمان ثبت خزیدن ذخیره می کند.

¹ - Crawler

حضور برچسب‌هایی که در گذشته‌ای دور بوده‌اند، در این مرحله، خوشه‌بندی کمرنگ شود.

$$\begin{aligned} \text{Sim}(Tag_i, Tag_j) & \\ = \frac{\beta \times \text{Sim}(Tag_i, Tag_j)}{\alpha + \beta} & \quad (3) \end{aligned}$$

, $\forall (Tag_i, Tag_j) \in M, i < j$

$$\begin{aligned} \text{Sim}(Tag_i, Tag_j) &= \text{Sim}(Tag_i, Tag_j) \\ &+ \frac{\alpha \times \text{Sim}_{new}(Tag_u, Tag_v)}{\alpha + \beta} \\ \forall (Tag_i, Tag_j) \in M, \forall (Tag_u, Tag_v) & \quad (5) \\ &\in M_{new} \\ \text{and } Tag_i == Tag_u, Tag_j == Tag_v, & \\ i < j, \quad u < v & \end{aligned}$$

که در آن:

M_{new} ؛ ماتریس مشابهت جدید به دست آمده توسط معادله (۴) است. M ؛ ماتریس مشابهت اولیه است.

Tag_i و Tag_j ؛ زوج برچسب‌های موجود در ماتریس M است. Tag_u و Tag_v ؛ زوج برچسب‌های موجود در پیمایش جدید هستند.

اگر چنانچه زوج برچسب جدیدی نسبت به ماتریس اولیه پدید آمده بود با توجه به معادله (۶) آن را به ماتریس اولیه اضافه می‌کنیم.

$$\begin{aligned} \text{Sim}(Tag_i, Tag_j) &= \frac{\alpha \times \text{Sim}_{new}(Tag_u, Tag_v)}{\alpha + \beta} \\ \forall (Tag_i, Tag_j) \notin M, \forall (Tag_u, Tag_v) & \quad (6) \\ &\in M_{new} \\ &, u < v \end{aligned}$$

۵. پس از تغییر ماتریس مشابهت، گراف مشابهت را تشکیل داده و با پارتیشن‌بندی مجدد آن با توجه به حد آستانه τ ، عمل خوشه‌بندی را برای تشکیل و یا حذف خوشه‌ها بررسی می‌کنیم که این شیوه را پارتیشن‌بندی گراف زمانی می‌نامیم.

۶. این عملیات دائماً می‌تواند در حال اجرا باشد که دلیل آن هم توقف نیافتن رشد اطلاعات جدید در وب است.

که در معادله فوق، M ؛ ماتریس مشابهت اصلی، α ؛ ضریب به‌روزرسانی و β ؛ ضریب به یادآوری گذشته است. ۲. عملیات خزنده را مجدداً از سر گرفته و مراحل آن را تکرار می‌کنیم تا این دفعه برچسب‌های جدیدی به دست آیند. اگر تعداد برچسب‌های جدید از حد عددی مانند بیشتر بودند به مرحله بعد می‌رویم؛ وگرنه به مرحله ۵ می‌رویم.

۳. در این مرحله - مشابه مازول خوشه‌بندی اولیه - K تا از پرکاربردترین برچسب‌ها که فراوانی آنها از حد عددی بیشتر باشد را انتخاب کرده و درجه شباهت آنها را مطابق معادله (۴) محاسبه می‌کنیم:

$$\text{Sim}_{new}(Tag_u, Tag_v) = \frac{|A_t \cap B_t|}{|A_t \cup B_t|} \quad (4)$$

که در آن A_t ؛ مجموعه سندهای جدید کشف‌شده در زمان t است که دارای Tag_u و به همین ترتیب B_t ؛ مجموعه سندهای جدید کشف‌شده در زمان t است که دارای Tag_v می‌باشند.

۴. مقادیر درایه‌هایی ماتریس جدید را با ماتریس اصلی مطابقت می‌دهیم. اگر زوج برچسبی هم در ماتریس اصلی حضور داشت و هم در ماتریس جدید، آنگاه درایه‌ای ماتریس اصلی را مطابق معادله (۵) به‌روز می‌کنیم.

۴. نتایج و ارزیابی

۴-۱- انتخاب مجموعه داده و استخراج

برچسب‌ها

مجموعه داده ما برگرفته از سایت MetaFilter [۱۳] می-باشد. این مجموعه داده شامل شماره لینک به همراه برچسب‌های ثبت شده برای آن توسط کاربر و همچنین زمان ثبت آن که بین بازه سال ۲۰۰۵ تا ۲۰۱۲ میلادی می‌باشد، است. در جدول (۱) جزئیات این مجموعه داده شرح داده شده است و در جدول (۲) تعدادی از پرکاربردترین برچسب‌ها آورده شده است.

جدول (۱): جزئیات مجموعه داده MetaFilter

| | |
|--|--------|
| تعداد لینک‌ها | ۱۰۲۹۴۵ |
| تعداد برچسب‌ها | ۱۲۰۷۰۵ |
| تعداد برچسب‌های نسبت داده شده به لینک‌ها | ۵۷۵۱۴۱ |

جدول (۲): تعدادی از پرکاربردترین برچسب‌ها در مجموعه داده

| MetaFilter | | | |
|------------|--------------------------|-----------------------|------------------------|
| برچسب | تعداد لینک نسبت داده شده | تعداد ماه میلادی مورد | تعداد روز مورد استفاده |
| brokenlink | 8450 | 27 | 255 |
| art | 4810 | 85 | 1934 |
| history | 2741 | 85 | 1414 |
| video | 2427 | 85 | 1450 |
| film | 1933 | 85 | 1141 |
| music | 5844 | 85 | 2054 |
| War | 1836 | 85 | 913 |
| politics | 3618 | 85 | 1390 |
| iraq | 1718 | 80 | 776 |
| internet | 1461 | 85 | 708 |

۴-۲- پیش پردازش اولیه

پس از بررسی اولیه داده‌ها متوجه شدیم که یکی از برچسب‌ها با وجود پر تکرارترین برچسب، ولی گویا مفهوم و استفاده خاصی ندارد. اگر به جدول (۲) مجدداً نگاهی بیندازید متوجه خواهید شد که برچسب "brokenlink" که به معنی لینک خراب یا شکسته است، ارتباطی با سایر

برچسب‌ها ندارد و خود سایت اصلی در صورتی که لینک مورد نظر خراب بوده، این برچسب را قرار داده است. لذا ما نیز از این برچسب در تمامی مراحل عملیات خوشه‌بندی صرف‌نظر می‌کنیم.

۴-۳- خوشه‌بندی اولیه از طریق پارتیشن‌بندی

گراف مشابهت

همان‌طور که بیان شد، خوشه‌بندی ما قرار است مبتنی بر تغییرات زمان عمل کند؛ لذا ابتدا خوشه‌بندی خود را به دو دوره تقسیم می‌کنیم:

۱. دوره اول مجموعه داده تا انتهای سال ۲۰۰۷،

۲. دوره ۲ تا ۵۰ مجموعه داده از ابتدای سال ۲۰۰۸ تا ماه

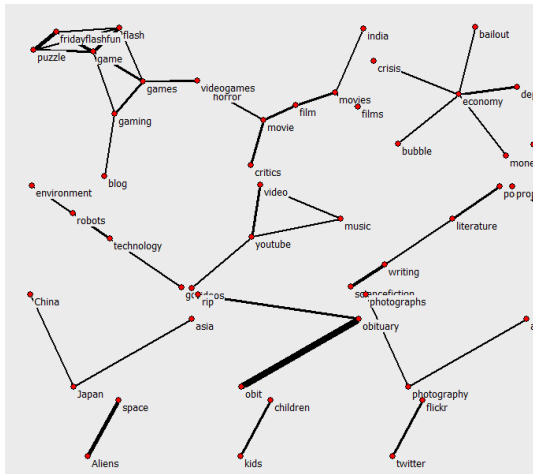
اول ۲۰۱۲ به تفکیک هر ماه یک دوره.

عملیات اولیه خوشه‌بندی ما بر روی مجموعه داده اولی یعنی ماقبل ۲۰۰۸ خواهد بود و برای خوشه‌بندی زمانی از ترکیب مجموعه اول و دوم و با حرکت بازه‌های زمانی ماه به ماه استفاده می‌کنیم. منطق این کار بدین دلیل است که فرض می‌کنیم یک مجموعه داده‌ای به ما داده شده است و حال از این به بعد قرار است با طی زمان، خوشه‌بندی را انجام دهیم. برای شروع کار ما $K_I=200$ برچسب از پرکاربردترین برچسب‌های دوره اول را انتخاب می‌کنیم. در این مجموعه داده، برخی مشخصات درجه تشابه به دست آمده برای ۲۰۰ برچسب مورد نظر مطابق جدول (۳) است.

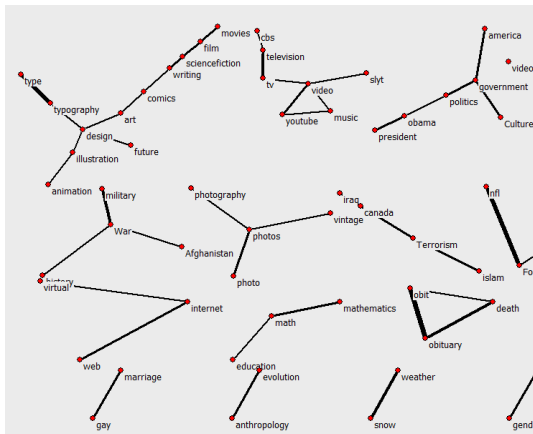
جدول (۳): برخی مشخصات درجه تشابه به دست آمده برای کل

مجموعه داده

| | |
|-----------|--------|
| Min: | 0.0001 |
| Max: | 0.3326 |
| Mean: | 0.0046 |
| Variance: | 0.0001 |
| Std.: | 0.0107 |



شکل (۶): گراف جدید برای دوره ۱۴ (ماه اول سال ۲۰۰۹) با حد آستانه <0.07 (با حفظ ارتباط با دوره های گذشته و جدید)



شکل (۷): گراف جدید برای دوره ۲۶ (ماه اول سال ۲۰۱۰) با حد آستانه <0.07 (با حفظ ارتباط با دوره های گذشته و جدید)

یعنی مثلاً اگر تا سال ۲۰۰۸ واژه war با Iraq و usa در ارتباط بوده، می بینیم که این ارتباط در شروع سال ۲۰۱۲ تغییر کرده و به جای Iraq، Iran اضافه شده است. این به منزله این است که کاربران با توجه به مسائل روز لینک ها و مطالبی که در سایت ها قرار می دهند، به روز می شود و مطابق با همین وضعیت برچسب های آنها نیز شکل دیگری به خود می گیرد. همان طور که اشاره شد، هدف ما در این مقاله نیز به روز بودن خوشه بندی در طی زمان با برچسب های درج شده است.

برای نیل به هدف به روز نگه داشتن خوشه ها، همان طور که گفتیم، مجموعه داده به دو بخش تقسیم شده بود. ابتدا ما خوشه بندی را بر روی مجموعه داده دوره اول (از سال ۲۰۰۵ تا انتهای ۲۰۰۷) انجام می دهیم. سپس برای هر بازه (ما بازه وجود ندارد) مطابق معادله های (۳) تا (۶) ماتریس مشابهت اصلی را به روز رسانی می کنیم.

جدول (۵): مقادیر پارامترهای روش پیشنهادی

| ماه | ۰/۰۵ | ۲۰۰ | ۲۰۰ | K_1 | K_2 |
|-----|------|-----|-----|-------|-------|
| ۱ | ۲ | ۲۰۰ | ۲۰۰ | | |

حال کافی است که از روی این ماتریس در هر بازه گراف را تشکیل داده تا پس از پار تیشن بندی آن خوشه بندی جدید حاصل شود. طبیعتاً انتظار ما از خوشه های به دست آمده این طور خواهد بود که ضمن حفظ ارتباطات مفاهیم گذشته، با مفاهیم روز نیز شدیداً در ارتباط باشد. ما برای خوشه بندی در اینجا نیز از Maximal Clique استفاده می کنیم. همچنین در جدول (۵) مقادیر پارامترهای روش پیشنهادی ذکر شده است.

برای مقایسه عملکرد، ما روش خود را Temporal Co-occurrence نامیده و با دو روش ذیل مقایسه می کنیم:

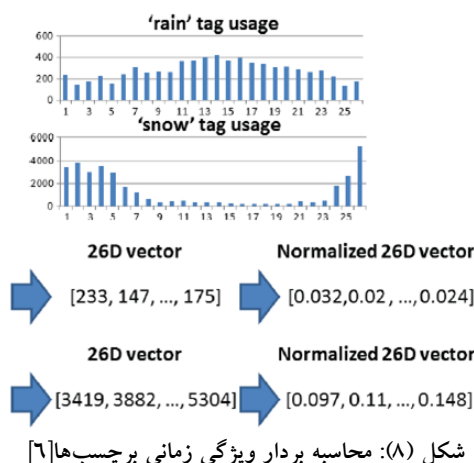
۱. روش Co-Occurrence در [3, 4, 5, 6, 7 و 8]

۲. روش Temporal Vector در [6]

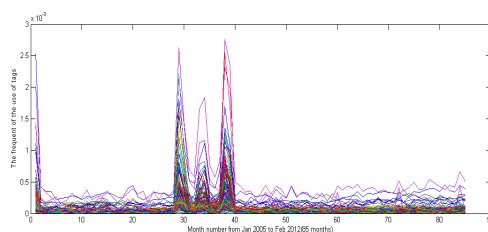
روش اول ما خوشه بندی را به این طریق انجام خواهیم داد که در هر دوره زمانی T ، ۲۰۰ برچسب پر کاربرد را پیدا کرده و برای آنها ماتریس مشابهت را تشکیل دهیم و پس از تشکیل گراف از روی ماتریس فوق و تعیین یک حد

آستانه ثابت برای تمام زمان‌ها، پاریشن‌بندی و در نتیجه خوشه‌بندی را انجام دهیم.

برای روش دوم برای هر دوره بردار ویژگی را برای ۲۰۰ برچسب پر کاربرد تا آن دوره را مطابق شکل (۸) تشکیل می‌دهیم. برای این کار تعداد تکرار هر برچسب را از اولین دوره تا آن دوره شمارش می‌کنیم و برای نرمال شدن، آن را بر مجموع تکرارهای کل برچسب‌های دوره تقسیم کرده و سپس برای خوشه‌بندی از روش $k=20$ استفاده می‌کنیم.



ما برای این بخش دوره‌های زمانی (بازه‌های زمانی) را ماهیانه و از اولین ماه سال ۲۰۰۵ تا اولین ماه سال ۲۰۱۲ در نظر گرفتیم که مجموعاً ۸۵ دوره شده است. در شکل (۹) توزیع ماهانه این برچسب‌ها قابل مشاهده است.



شکل (۹): توزیع ۲۰۰ برچسب از پرکاربردترین برچسب‌ها از ابتدای ۲۰۰۵ تا ابتدای ۲۰۱۲ به تفکیک ماه

پس از انجام خوشه‌بندی در سر دوره‌های زمانی مشخص شده با ۳ روش فوق‌ما، لازم است معیاری برای مقایسه این روش‌ها اتخاذ کنیم.

۵-۴- تعیین معیار ارزیابی

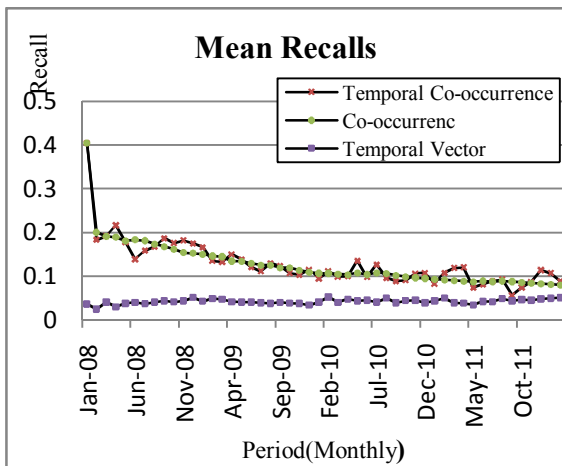
برای محاسبه عددی ارزیابی روش‌های مورد مقایسه، ما از معیارهای Recall و Precision استفاده می‌کنیم. معیار Precision به ما می‌گوید که چه کسری از تعداد نمونه‌های به دست آمده^۱ مورد نظر در خوشه، با یکدیگر مرتبط^۲ هستند. (صحت) و Recall یعنی چه کسری از تعداد نمونه‌های مرتبط واقعی در خوشه قرار دارد. همچنین از معیار F-Measure برای نمایش ارتباط بین Recall و Precision [۵ و ۱۵] نیز استفاده خواهیم کرد که معادله آنها به طریق ذیل است:

که در آن R ، مجموعه برچسب‌های مرتبط و G ؛ مجموعه برچسب‌های موجود در هر خوشه است.

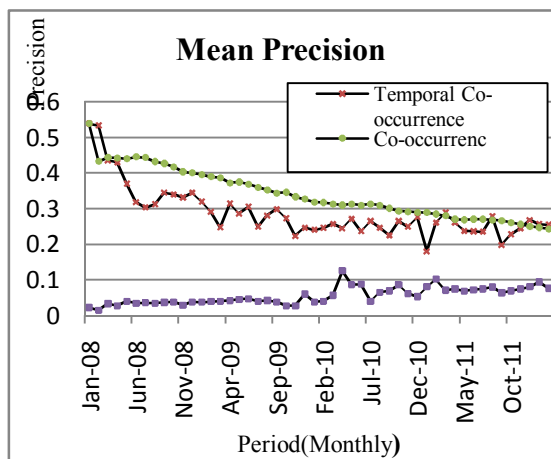
لازم به ذکر است در معیارهای Recall و Precision نکته قابل تأمل نحوه محاسبه آنهاست. در [۶] برای تصمیم‌گیری برچسب‌های مرتبط از طریق طراحی یک سایت، از مردم نظرسنجی کرده‌اند، ولی به علت محدودیت این شیوه و یا عدم اطمینان از صحت درستی نظرات مردم، ما به دو روش آنها را برای هر خوشه محاسبه می‌کنیم. سپس در بازه زمانی مورد نظر، میانگین معیارهای محاسبه شده خوشه‌های هر روش خوشه‌بندی را به دست می‌آوریم و با یکدیگر مقایسه می‌کنیم.

¹- Retrieved instances
²- Relevant

با توجه به این محاسبه، در شکل (۱۰) مقایسه میانگین Recall خوشه‌ها برای هر سه روش گفته شده و در شکل (۱۱) مقایسه میانگین Precision خوشه‌ها برای همان سه روش آورده شده است. در جدول (۶) نیز مقایسه میانگین مقادیر Precision، Recall و F-Measure برای تمام خوشه‌ها در تمامی بازه‌های زمانی جهت مقایسه روش پیشنهادی آورده شده است.



شکل (۱۰): مقایسه Mean Recalls خوشه‌ها در بازه‌های زمانی مختلف (روش محاسبه Recall مستقل از زمان)



شکل (۱۱): مقایسه Mean Precisions خوشه‌ها در بازه‌های زمانی مختلف (روش محاسبه Precision مستقل از زمان)

$$\text{Precision} = \frac{|R \cap G|}{|R|} \quad (7)$$

$$\text{Recall} = \frac{|R \cap G|}{|G|} \quad (8)$$

$$F = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (9)$$

برای هر دو روش محاسبه Precision و Recall در هر

خوشه داریم:

t : بازه‌ی زمانی (شماره ماهی که از ماه اول سال 2008

گذشته است) که $1 \leq t \leq 50$

Tag_i : هر برچسب موجود در خوشه،

G : تمام برچسب‌های خوشه به غیر از Tag_i

A : مجموعه محاسبه شده ماتریس‌های مشابهت برای هر

دوره زمانی که درجه مشابهت آنها از حد آستانه ذکر شده

بیشتر باشد و دارای ۴ فیلد است. برچسب اول، برچسب دوم،

درجه تشابه، عدد بازه زمانی.

B : زیرمجموعه‌ای از A که دارای شرایط ذیل باشد:

- عدد بازه زمانی آن از t کمتر باشد.

- برچسب اول با Tag_i برابر باشد.

۱-۵-۴- روش اول محاسبه Precision و Recall

مستقل از زمان

در این روش در مجموعه B بیشترین تکرار شماره بازه

زمانی را پیدا می‌کنیم (peak در بردار زمانی برای این

برچسب مطابق شکل (۹) - یعنی بیشترین برچسب‌هایی که

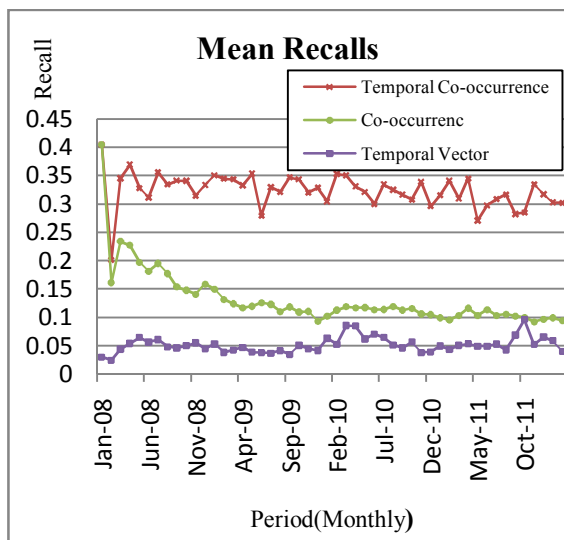
با برچسب Tag_i مستقل از زمان شباهت زیادی داشتند - و آن

را f می‌نامیم. سپس تعریف می‌کنیم:

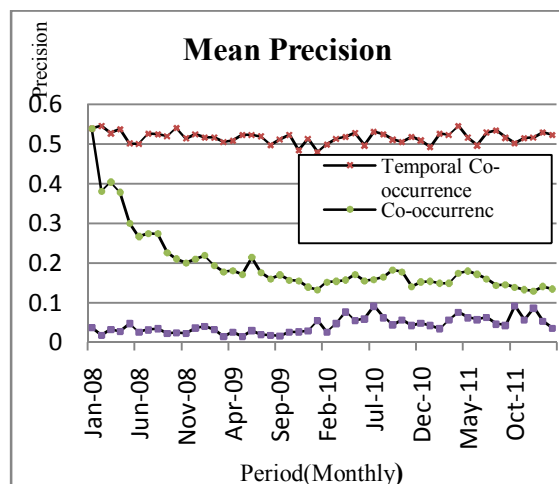
R : زیرمجموعه‌ای از برچسب‌های B که عدد بازه زمانی

آن برابر با f باشد.

$RelevantSize$: اندازه اشتراک مجموعه R با مجموعه G



شکل (۱۲): مقایسه Mean Recalls خوشه‌ها در بازه‌های زمانی مختلف (روش محاسبه Recall وابسته به زمان)



شکل (۱۳): مقایسه Mean Precisions خوشه‌ها در بازه‌های زمانی مختلف (روش محاسبه Precision وابسته به آخرین زمان)

جدول (۷): مقایسه معیارهای Precision, Recall و F- Measure (روش محاسبه Recall و Precision وابسته به زمان)

| Method | Recall | Precision | F-Measure |
|------------------------|----------|-----------|-----------|
| Temporal Co-occurrence | 0.323629 | 0.516722 | 0.397991 |
| Co-occurrence | 0.131678 | 0.195546 | 0.157379 |
| Temporal Vector | 0.051034 | 0.042432 | 0.046337 |

جدول (۶): مقایسه معیارهای Precision, Recall و F- Measure (روش محاسبه Recall و Precision مستقل از زمان)

| Method | Recall | Precision | F-Measure |
|------------------------|----------|-----------|-----------|
| Temporal Co-occurrence | 0.122651 | 0.287737 | 0.17199 |
| Co-occurrence | 0.121229 | 0.340581 | 0.178811 |
| Temporal Vector | 0.041395 | 0.055401 | 0.047385 |

۲-۵-۴- روش دوم محاسبه Precision و Recall وابسته به آخرین زمان

در این روش در مجموعه B به جای بیشترین تکرار شماره بازه زمانی، بزرگ‌ترین عدد بازه زمانی را پیدا می‌کنیم - یعنی اخیرترین برچسب‌هایی که با برچسب Tag شباهت زیادی داشتند - و آن را می‌نامیم. سپس تعریف می‌کنیم:

R : زیرمجموعه‌ای از برچسب‌های B که عدد بازه زمانی آن برابر با f باشد.

$RelevantSize$: اندازه اشتراک مجموعه R با مجموعه G مشابه نحوه محاسبه اول. در اینجا هم ضمن محاسبه معیارها با این روش، در شکل (۱۲) مقایسه میانگین Recall خوشه‌ها و در شکل (۱۳) مقایسه میانگین Precision خوشه‌ها قابل مشاهده است. در جدول (۷) نیز مقایسه میانگین مقادیر Recall، Precision و F-Measure برای تمام خوشه‌ها آورده شده است.

به ازای هر دو روش محاسبه، شکل‌ها و جداول آورده شده نشان دهنده مقایسه بین روش پیشنهادی با سایر روش‌هاست.

همان‌طور که در مقایسه شکل (۱۰) با شکل (۱۲) و همچنین مقایسه شکل (۱۱) با شکل (۱۳) مشخص است، می‌بینیم که در محاسبه معیارهای Recall و Precision به روش مستقل از زمان، به غیر از بازه‌هایی، زیاد شاهد بهبود نیستیم و علت آن هم بدین دلیل است که ما از دو ضریب α و β استفاده کردیم و به تدریج تأثیر گذشته را کم‌رنگ‌تر کردیم؛ حتی اگر در گذشته دور خیلی با هم ارتباط داشتند. البته اگر $\alpha = 1$ و $\beta = 0$ انتخاب می‌شد، نتیجه با روش Co-occurrence برابر می‌شد؛ بنابراین اگر چه ممکن است دو برچسب در یک بازه زمانی به اوج شباهت خود رسیده باشند، ولی در آینده تضمینی برای حفظ این ارتباط وجود ندارد. درحالی‌که در محاسبه معیارهای Recall و Precision به روش وابسته به زمان می‌بینیم که روش پیشنهادی نسبت به دو روش دیگر نتیجه بهتری داشته است.

نکته دیگری که می‌توان استنتاج کرد این است که در روش Co-occurrence در ابتدا نتایجی خوبی را از خود نشان می‌دهد، ولی به تدریج رو به کاهش خواهد بود؛ بنابراین اگر چندین دوره تناوب زمانی دیگر روش‌ها را تکرار می‌کردیم مطمئناً در روش محاسبه معیارهای Recall و Precision به روش مستقل از زمان، روش پیشنهادی بهتر از Co-occurrence می‌شد؛ هر چند انتهای همین بازه‌های مشخص شده، گویای این ادعاست. در نهایت با توجه به ارزیابی‌های انجام شده شاهد آن هستیم که نتیجه کار پیشنهادی به علت تطبیق‌پذیری با زمان نتایج مطلوب‌تری نسبت به سایر روش‌ها از خود به‌روز داده است.

۵- نتیجه‌گیری و پیشنهادها

در عصر حاضر با ظهور وب ۲.۰ ما شاهد رشد چشمگیری فعالیت‌های کاربران در سطح وب هستیم. یکی از عمده‌ترین این فعالیت‌ها به اشتراک گذاشتن منابع با سایر کاربران است. در بعضی از خدمات وب ۲.۰ از فکسونومی برای رساندن مفهوم بهتر و خلاصه‌تر، یعنی امکان برچسب‌گذاری منابع توسط کاربران قرار داده شده است که در چند سال اخیر نیز استقبال خوبی از آن شده است؛ لذا با توجه به حجم وسیع منابع، این برچسب‌ها ملاک بسیار خوبی برای خوشه‌بندی و دسته‌بندی جهت دسترسی کاربران به منابع مرتبط می‌باشد. از طرف دیگر ممکن است روزانه برچسب‌های جدیدی به بانک برچسب‌ها اضافه شود و یا حتی ارتباطات برچسب‌ها در طی زمان تغییر کنند؛ بنابراین نیاز به روش (هایی) است که خوشه‌بندی‌ها را با توجه به تغییرات زمانی به‌روز نگه دارد.

ما در این مقاله روش‌های مختلف خوشه‌بندی برای برچسب‌ها را در سیستم‌های برچسب‌گذاری مستقل از زمان و وابسته به زمان، شرح دادیم. سپس روش پیشنهادی خود که وابستگی برچسب‌ها نسبت به زمان را خوشه‌بندی می‌کرد، شرح دادیم و در پایان ضمن مقایسه با سایر روش‌های مشابه شاهد بهبود کیفی و کمی بهتر خوشه‌های به وجود آمده بودیم و دیدیم که روش پیشنهادی، F-Measure میانگین خوشه‌ها را در حالت وابسته به زمان نسبت به بهترین روش ۲۴٪ بهبود بخشید. در پایان، پیشنهاد ما برای کارهای آینده استفاده از ترکیب کلمات کلیدی استخراج شده از منابعی است که متنی هستند و وزن دهی مناسب به آنها و همچنین استفاده از آنتولوژی WordNet برای استخراج ریشه کلمات کلیدی استخراج شده و اضافه کردن آنها به مجموعه برچسب‌ها برای افزایش صحت خوشه‌ها می‌باشد. همچنین پیشنهاد استفاده از نتایج خوشه-

های به دست آمده از این روش را برای سیستم‌های
پیشنهاددهنده و دسته‌بندی منابعی که برای آنها

برچسب‌گذاری انجام شده، داریم.

۶- مراجع

- 1- Svein Anfinnsen, Gheorghita Ghinea, Sergio de Cesare, “**Web 2.0 and folksonomies in a library context**” *International Journal of Information Management* 31 (2011) 63–70
- 2- O’Reilly, T. (2006). “**Web 2.0 compact definition: Trying again**”. Available at: http://radar.oreilly.com/archives/2006/12/web_20_compact.html.
- 3- Golder, S., and Huberman, B. 2006. “**Usage patterns of collaborative tagging systems**”. *Journal of Information Science* 32(2):198–208.
- 4- <http://kmttoolbox.ir/?p=58>
- 5- Nitin Agarwal, Magdiel Galan, H. L., and Subramanya, S. “**WisColl: Collective wisdom based blog clustering**”. *Information Sciences* 180 (2010) 39–61, journal homepage: www.elsevier.com/locate/ins
- 6- Haipeng Zhang, Mohammed Korayem, Erkang You, David J. Crandall, “**Beyond co-occurrence: discovering and visualizing tag relationships from geo-spatial and temporal similarities**”, *WSDM '12 Proceedings of the fifth ACM international conference on Web search and data mining* Pages 33-42-2012
- 7- Micah Dubinko, Ravi Kumar, Joseph Magnani, Jasmine Novak, Prabhakar Raghavan, Andrew Tomkins, “**Visualizing tags over time**” *WWW '06 Proceedings of the 15th international conference on World Wide Web*. Pages 193 - 202, 2006
- 8- Edwin Simpson, HP Laboratories, Bristol, “**Clustering Tags in Enterprise and Web Folksonomies**”, *International Conference on Weblogs & Social Media, Seattle, March 31st, 2008*
- 9- Grigory Begelman, Philipp Keller, Frank Smadja, “**Automated Tag Clustering: Improving search and exploration in the tag space**” *ICWE'11 Proceedings of the 11th international conference on Web engineering*
- 10- Alberto Pérez García-Plaza a, Arkaitz Zubiaga, Victor Fresno, Raquel Martínez, “**Reorganizing clouds: A study on tag clustering and evaluation**”, *Journal of Expert Systems with Applications* 39 (2012) 9483–9493
- 11- Brooks, C. H., and Montanez, N. “**Improved annotation of the blogosphere via autotagging and hierarchical clustering**”. *Proceedings of the 15th international conference on World Wide Web, New York, USA (2006)*, pp. 625–632.
- 12- Jaiswal, Mayank Prakash, “**Clustering Blog Information**” (2007). *Master's Projects. Paper 36*, http://scholarworks.sjsu.edu/etd_projects/36
- 13- <http://stuff.metafilter.com/infodump/>
- 14- <http://pajek.imfm.si/>
- 15- http://en.wikipedia.org/wiki/Precision_and_recall