

تحلیل محتوا، رویکردی نوین در بهبود کارایی تشخیص اجتماع

علی ربیعانی^۱، حسین علیزاده^۲، بهروز مینایی^۳

^۱ دانشکده فناوری اطلاعات، دانشگاه علوم و فنون مازندران، بابل، areihanian@ustmb.ac.ir
^۲ دانشکده مهندسی کامپیوتر، دانشگاه علم و صنعت ایران، تهران، halizadeh@iust.ac.ir
^۳ دانشکده مهندسی کامپیوتر، دانشگاه علم و صنعت ایران، تهران، b_minai@iust.ac.ir

جدول زیر نیز مشخصات شبکه‌های ایجاد شده در هر یک از چارچوب‌ها را نشان می‌دهد:

جدول (۱): مشخصات شبکه‌های ایجاد شده

تعداد یال‌ها	تعداد گره‌ها	نام خوشه	چارچوب
۷۶۷۰۵	۵۹۲	-	چارچوب ۱
۱۵۸۳۳	۳۵۲	مستند	چارچوب ۲
۶۹۲۶۹	۴۹۱	وسترن	چارچوب ۳
۷۶۷۰۵	۵۹۲	-	چارچوب ۴

چارچوب اول، بیانگر حالتی است که ما تنها اقدام به اعمال الگوریتم تشخیص اجتماع بر روی شبکه‌ی اصلی می‌کنیم و هیچگونه تحلیل محتوایی صورت نمی‌پذیرد. چارچوب دوم همان تشخیص اجتماع مبتنی بر موضوع بوده و چارچوب سوم نیز تشخیص اجتماع بر روی شبکه معنایی می‌باشد از آنجایی که در چارچوب دوم، ما اقدام به یافتن خوشه‌های موضوعی می‌کنیم، بنابراین پس از نسبت دادن گره‌ها به خوشه‌های موضوعی، شبکه‌ی اصلی به دو بخش تقسیم می‌شود که هر یک از این بخش‌ها تنها شامل گره‌هایی خواهند بود که دارای موضوع مشترکند. لذا در جدول بالا در قسمت مربوط به چارچوب دوم، شبکه به دو بخش مستند و وسترن تقسیم شده است.

در نهایت، با توجه به رابطه‌ی (۱)، میزان $PurQ_{\beta}$ را در هر مرحله محاسبه می‌کنیم. جدول زیر، نشان‌دهنده‌ی نتایج نهایی اعمال چارچوب‌های مختلف بر روی مجموعه داده‌ی MovieLens می‌باشد:

جدول (۲): نتایج نهایی اعمال چارچوب‌های مختلف

مقدار $PurQ_{\beta}$	مقدار $Purity$	مقدار بودمی کل	مقدار بودمی خوشه	نام خوشه	چارچوب
۰.۱۹۱۵۳۵	۰.۹۷۴۵۲۴	۰.۱۰۸۶	۰.۱۰۸۶	-	چارچوب ۱
۰.۲۲۱۳۲۲	۱	۰.۱۲۴۴۲۷	۰.۲۶۸۶	مستند	چارچوب ۲
۰.۲۵۰۴۲۲	۰.۸۹۴۷۹۸	۰.۱۶۲۷	۰.۱۶۲۷	وسترن	چارچوب ۳

از آنجایی که ما به دنبال اجتماعاتی هستیم که این اجتماعات دربرگیرنده‌ی اعضای باشند که این اعضا هم دارای ارتباط تنگاتنگ با یکدیگر و هم دارای موضوع مشترک باشند، لذا میزان خلوص موضوعات و ساختار شبکه را ما به یک اندازه اهمیت داریم. بنابراین، برای محاسبه‌ی $PurQ_{\beta}$ مقدار β را برابر ۱ در نظر گرفتیم.

نتیجه‌ای که می‌توان از جدول بالا بدست آورد، این است که در هر دو چارچوبی که به تحلیل محتوا پرداختند نسبت به حالتی که تحلیل محتوایی انجام نشد، میزان بودمانی افزایش داشته است، با این تفاوت که در چارچوب سوم، با توجه به تحلیل محتوای مربوط به هر گره و اعمال مستقیم آن به صورت زوددهی بر روی یال‌ها، تاثیر بیشتری را بر روی مقدار بودمانی شاهد بودیم. همچنین، واضح است که در چارچوب دوم، از آنجایی که خوشه‌های موضوعی جدا شدند و الگوریتم تشخیص اجتماع بر روی این خوشه‌ها اعمال گردید، میزان خلوص موضوعات ($Purity$) بالاتر از دو چارچوب دیگر می‌باشد. در مجموع، با در نظر گرفتن میزان $PurQ_{\beta}$ برای ۳ چارچوب پیاپی مشاهده شد، واضح است که چارچوب‌هایی که تنها به ساختار و پیوند توجه نکرده‌اند و به تحلیل محتوا نیز پرداخته‌اند (چارچوب‌های ۲ و ۳)، به نتایج بهتری منجر شده‌اند. بنابراین می‌توان نتیجه گرفت که استفاده از نتایج تحلیل محتوا، موجب بهبود کارایی عمل تشخیص اجتماع می‌گردد.

نتیجه‌گیری:

در این مقاله، به بررسی روش‌هایی که از نتایج حاصل از تحلیل محتوای یک شبکه به منظور بهبود تشخیص اجتماعات استفاده کردند، پرداخته شد. بدین منظور، دو چارچوب تشخیص اجتماع مبتنی بر عنوان و تشخیص اجتماع بر روی شبکه معنایی پیاده‌سازی شده و بر روی مجموعه داده‌ی MovieLens اعمال شدند. نتایج حاصل از اعمال این دو چارچوب و مقایسه‌ی آن با نتایج بدست آمده از حالتی که تشخیص اجتماع بدون تحلیل محتوا صورت می‌پذیرد، بیانگر این واقعیت بود که تحلیل محتوا در یک شبکه و استفاده از نتایج آن منجر به بهبود عمل تشخیص اجتماع می‌گردد.

از آنجایی که تشخیص اجتماع مبتنی بر عنوان، نسبت به تشخیص اجتماع بر روی شبکه معنایی و نسبت به روشی که تشخیص اجتماع، بدون تحلیل محتوای شبکه صورت می‌پذیرد، منجر به یافتن اجتماعاتی با میزان خلوص موضوعات ($Purity$) بالاتری شد و از طرف دیگر تشخیص اجتماع بر روی شبکه معنایی، نسبت به تشخیص اجتماع مبتنی بر عنوان و نسبت به روشی که تشخیص اجتماع، بدون تحلیل محتوای شبکه صورت می‌پذیرد، منجر به رسیدن به مقدار بودمانی بالاتری شد، لذا در کارهای آینده بنا داریم تا چارچوبی ارائه دهیم که از مزیت مربوط به تشخیص اجتماع مبتنی بر عنوان و تشخیص اجتماع بر روی شبکه معنایی به طور همزمان استفاده کند. یعنی این چارچوب بتواند به طور همزمان هم میزان بودمانی و هم میزان خلوص موضوعات را افزایش دهد.

منابع:

- Leskovec, J., K.J. Lang, and M. Mahoney. "Empirical comparison of algorithms for network community detection", in Proceedings of the 19th international conference on World Wide Web. 2010. ACM.
- Zhao, Z., et al., "Topic oriented community detection through social objects and link analysis in social networks", Knowledge-Based Systems, 2012. 26: p. 164-173.
- Xia, Z. and Z. Bu, "Community detection based on a semantic network". Knowledge-Based Systems, 2012. 26: p. 30-39.
- Easley, D. and J. Kleinberg, *Networks, crowds, and markets*, Vol. 8. 2010: Cambridge Univ Press.
- Newman, M., *Communities, "modules and large-scale structure in networks"*, Nature Physics, 2011.
- Porter, M.A., J.-P. Onnela, and P.J. Mucha, "Communities in networks", Notices of the AMS, 2009. 56(9): p. 1082-1097.
- Lancichinetti, A. and S. Fortunato, "Consensus clustering in complex networks", Scientific reports, 2012. 2.

چکیده:

یکی از مباحث مهم در زمینه‌ی تحلیل شبکه‌های پیچیده، بحث تشخیص اجتماع می‌باشد. اکثر روش‌هایی که در زمینه‌ی تشخیص اجتماع پیشنهاد شده‌اند، این عمل را تنها با در نظر گرفتن ساختار گراف شبکه انجام می‌دهند، اما در سال‌های اخیر، تلاش‌هایی در زمینه‌ی تحلیل محتوا و استفاده از نتایج آن به منظور بهبود کارایی تشخیص اجتماع صورت گرفته است. در این مقاله برآنیم تا با معرفی برخی از این رویکردهای جدید و پیاده سازی روش پیشنهادیمان، به بررسی نتایج حاصل از بکارگیری تحلیل محتوا در عمل تشخیص اجتماع بپردازیم. نتایج حاصله نشان می‌دهند که بکارگیری محتوا به صورت ملموسی به بهبود کارایی عمل تشخیص اجتماع کمک می‌کند.

کلمات کلیدی

شبکه پیچیده، اجتماع، تشخیص اجتماع، تحلیل محتوا

مبانی نظری تحقیق:

نحوه نمایش شبکه: می‌توان شبکه‌ها را به دو صورت نمایش داد، به صورت گراف و به صورت ماتریس. در نمایش گرافی شبکه، ساختارهای گرافی، وظیفه متصل کردن گره‌ها و یال‌ها به یکدیگر را بر عهده دارند. به منظور نمایش ماتریسی شبکه، یک ماتریس مربعی ایجاد می‌شود که به آن ماتریس مجاورتی گراف شبکه گفته می‌شود. ابعاد این ماتریس، برابر تعداد گره‌های گراف مشناظر با آن بوده و در صورت وجود یال بین دو گره، در درایه‌ی متناظر با دو گره در ماتریس مجاورتی، مقدار ۱ لحاظ شده و در غیر اینصورت در این درایه، مقدار ۰ لحاظ می‌شود.

تعریف اجتماع: در مقالات و تحقیقات علمی مختلف، تعاریف متفاوتی از اجتماع ارائه شده است. یکی از این تعاریف که در مقاله‌ی ۱ آمده است بیان می‌کند که یک اجتماع شبکه که گاهی از اوقات به آن پیمان یا خوشه گفته می‌شود، به عنوان یک گروه از گره‌ها در نظر گرفته می‌شود که تعاملات بهتر و بیشتری بین اعضای این گروه نسبت به اعضای این گروه و دیگر اعضای شبکه وجود دارد.

تابع بودمانی: مشهورترین معیار ارزیابی اجتماعات تشخیص داده شده توسط یک الگوریتم تشخیص اجتماع، بودمانی می‌باشد. پس از این‌که تمامی اجتماعات یک شبکه شناسایی شدند، بودمانی به آن شبکه اعمال می‌شود. ورودی تابع بودمانی، شبکه و تمامی اجتماعات آن بوده و خروجی آن، یک عدد حقیقی بین ۰ و ۱ می‌باشد. مقادیر بودمانی نزدیک به ۱ برای یک روش تشخیص اجتماع، بیانگر مناسب بودن آن روش می‌باشد. اگر تمامی رنوس در یک اجتماع قرار گرفته باشند (یعنی کل گراف شامل تنها یک اجتماع باشد) و یا گره‌ها بصورت تصادفی در بین اجتماعات قرار گرفته باشند، مقدار بودمانی برابر ۰ خواهد بود. اگر روش تشخیص اجتماع، هر راس گراف را در یک اجتماع مجزا قرار دهد، در صورتی که شبکه واقعا چنین ساختاری را نداشته باشد، آنگاه مقدار بودمانی نزدیک به ۱- خواهد بود.

تابع $PurQ_{\beta}$: در مقاله‌ی ۲، چارچوبی پیشنهاد شد که این چارچوب مباحث موضوع و پیوند را با یکدیگر ترکیب کرده و به عمل تشخیص اجتماع می‌پردازد. لذا به منظور ارزیابی نتایج حاصل از این چارچوب پیشنهادی، معیار جدید معرفی شد. بدیهی است که چنین معیاری، یاستی به بررسی اجتماعات از دو جنبه‌ی موضوع و پیوند بپردازد. این معیار که با نام $PurQ_{\beta}$ معرفی شده است، به صورت زیر تعریف می‌شود:

$$PurQ_{\beta} = (1 + \beta^2) \cdot (Purity \cdot Q) / (\beta^2 \cdot Purity + Q) \quad (1)$$

در رابطه‌ی بالا، $Purity$ ، میزان خلوص موضوعات، در اجتماعات تشخیص داده شده را مشخص می‌کند. هرچقدر مقدار $Purity$ بالاتر باشد، نشان‌دهنده‌ی این امر خواهد بود که اجتماعات، از جنبه‌ی موضوع، بهتر افراز بندی شده‌اند. Q نیز نمایانگر بودمانی می‌باشد. همچنین در رابطه‌ی بالا، β پارامتری برای تنظیم وزن $Purity$ و Q می‌باشد و می‌تواند دارای مقادیر از ۰ تا بی‌نهایت باشد. اگر مقدار β برابر با ۱ باشد، می‌تواند اینگونه در نظر گرفته شود که میزان خلوص موضوعات و ساختار شبکه اجتماعی به یک اندازه اهمیت دارند.

روش تحقیق:

پژوهش‌گران مختلف الگوریتم‌های گوناگونی را به منظور تشخیص اجتماع پیشنهاد داده‌اند که این الگوریتم‌ها، تنها با در نظر گرفتن ساختار گرافی شبکه به عمل تشخیص اجتماع می‌پردازند. در سال‌های اخیر، تلاش‌هایی در زمینه‌ی تحلیل محتوای یک شبکه و استفاده از نتایج آن، در بهبود عمل تشخیص اجتماع انجام گرفته است. به عنوان نمونه‌ای از این تلاش‌ها، می‌توان به رویکردی که تحت عنوان تشخیص اجتماع مبتنی بر موضوع مطرح شده است، اشاره کرد. در این رویکرد، ابتدا خوشه‌های موضوعی در شبکه یافت می‌شوند و سپس، یکی از الگوریتم‌های تشخیص اجتماع بر روی این خوشه‌های موضوعی اعمال می‌شود. به عنوان نمونه‌ی دیگری از این تلاش‌ها، می‌توان به رویکردی که تحت عنوان تشخیص اجتماع بر روی شبکه معنایی مطرح شده است، اشاره کرد. در این رویکرد، ابتدا محتوای رد و بدل شده بین عناصر موجود در شبکه، تحلیل و بررسی می‌شوند و بر اساس نتایج این تحلیل، یک شبکه معنایی ترسیم می‌شود و عمل تشخیص اجتماع بر روی این شبکه معنایی انجام می‌پذیرد. در این مقاله، این دو رویکرد به طور مبسوط بررسی شده و روش‌های پیشنهادیمان پیاده‌سازی شدند.

یافته‌های تحقیق:

در این تحقیق، چارچوب‌های ارائه شده در روش‌های تشخیص اجتماع مبتنی بر موضوع و تشخیص اجتماع بر روی شبکه معنایی، بر روی یکی از مجموعه داده‌های MovieLens اعمال شدند. این مجموعه داده، شامل ۱۰۰۰۰۰ امتیازدهی اعمال شده به ۱۶۸۲ فیلم توسط ۹۴۳ کاربر می‌باشد این امتیازها از ۱ تا ۵ می‌باشند. به منظور پیاده‌سازی چارچوب‌های ذکر شده، کلیه فیلم‌هایی که در موضوع مستند یا موضوع وسترن یا در تلفیقی از هر دو موضوع ساخته شده بودند، انتخاب شدند. در واقع در این مجموعه داده، فیلم‌ها به عنوان اشیاء اجتماعی در نظر گرفته شدند. در کل این مجموعه داده، ۷۷ فیلم در یکی از این دو موضوع (مستند و وسترن) و ۱۱۰ فیلم در هر دو موضوع، ساخته شده‌اند (برای مثال، امکان دارد فیلم مستندی به سبک وسترن ساخته شده باشد. در این حالت، این فیلم هر دو موضوع را شامل می‌شود). سپس کلیه امتیازهایی که توسط کاربرین به این فیلم‌ها اختصاص داده شده است، بازبینی شدند. لازم به ذکر است که احتمال دارد کاربری به بیش از یک فیلم امتیاز داده باشد. سپس شبکه‌ای رسم می‌شود که در آن، بین تمامی کاربرانی که به فیلم‌ها مشترکی به امتیاز دادند، یال وجود خواهد داشت. وزن این یال‌ها برابر تعداد فیلم‌های مشترکی خواهد بود که دو گره (کاربر) مربوطه به امتیازدهی به آن‌ها پرداختند. شکل زیر نمای از این شبکه (که ما به آن شبکه‌ی اصلی می‌گوییم) را نشان می‌دهد:



شکل (۱): نمای از شبکه‌ی اصلی