

مروری بر الگوریتمهای انجمن یابی در تحلیل شبکه های اجتماعی

بهروز ساعی^{۱*}، علیرضا نوروزی^۲

۱-دانشجوی کارشناسی ارشد فناوری اطلاعات، گروه کامپیوتر، دانشگاه آزاد اسلامی شهر مجلسی

Email: behrozsaei1@yahoo.com

۲-استادیار، گروه کامپیوتر، دانشگاه آزاد اسلامی واحد شهر مجلسی

Email: Norouzi.arz@gmail.com

چکیده:

داده کاوی یکی از پشرفتهای اخیر در حوزه کامپیوتر برای اکتشاف عمقی داده هاست. داده کاوی اطلاعات حیاتی و مهمی را که برای برنامه ریزی های استراتژیک، مورد نیاز است، آشکار میکند. برخی شبکه های اجتماعی (و اغلب شبکه های اجتماعی برخط) که خاصیت انجمنی (یعنی دارای انجمنهایی هستند) را از خود بروز میدهند، دارای ساختار سلسله مراتبی هستند. بدین معنا که انجمن های بزرگتر از یکمجموعه انجمن های کوچکتر تشکیل شده و این انجمن های کوچکتر خود از یک دسته انجمن های کوچکتر دیگر تشکیل شده است. شبکه های اجتماعی اغلب این خاصیت را از خود بروز می دهند. این مقاله سعی دارد تعدادی از الگوریتمهای مهم در حوزه انجمن کاوی را که در کاوش شبکه های اجتماعی، بکار می روند، معرفی کند. الگوریتمهای تقسیم، ماجولاریتی، شعاع طیفی، مقادیر ویژه، پویا، استنباط آماری و سایر روشها، در این مقاله به صورت کلی مرور و معرفی میشوند.

کلمات کلیدی: داده کاوی، شبکه های اجتماعی، الگوریتم های، تقسیم، ماجولاریتی، شعاع طیفی، مقادیر ویژه، پویا، استنباط آماری

۱-مقدمه

یکی از سخت ترین کارهای مربوط به انجمن یابی در شبکه اجتماعی در همان گام نخست روی می دهد؛ سوال اساسی و پایه ای مربوط به تعریف انجمن هاست. شکل ۱ تصویر یک زیرگراف از گراف فیس بوک است که از ۴۰۳۹ راس ۸۲۳۴ یال تشکیل شده می باشد. هرچقدر به لحاظ بصری یافتن انجمن های این گراف راحت به نظر می رسد اما به لحاظ علم ریاضی و نظریه گراف ارائه تعریف و یافتن این انجمن ها پیچیده می نماید. در گراف شکل ۱ توده های مجتمع راس ها، انجمن ها را تشکیل می دهند. شکل ۲ یک گراف تصادفی است با تعداد راس ها و یال های.

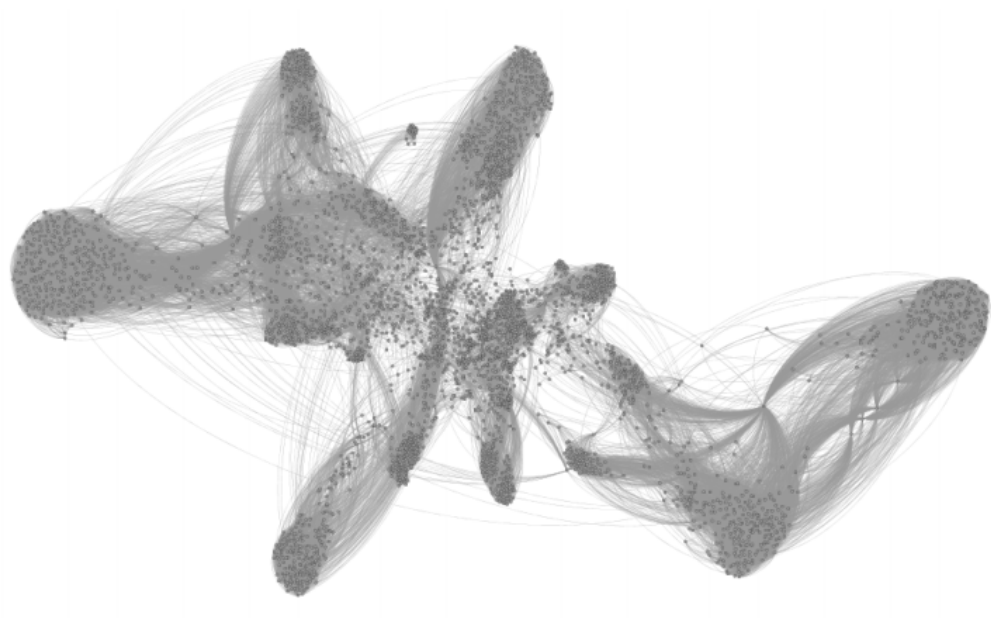
مساوی با گراف فیس بوک که برای هر یال دو راس را به طور تصادفی انتخاب کرده و به هم وصل می کند. همانطور که مشخص است. به لحاظ شهودی و تصویری گراف های مربوط به شبکه اجتماعی از گروه های متراکمزیادی تشکیل شده است اما گراف مربوط به گراف تصادفی تنها از یک توده تشکیل شده است. این دو گراف تفاوت های دیگری نیز دارند که در ادامه به آن خواهیم پرداخت، از جمله این تفاوت ها وجود چند راس (قطب) است که نقش ارتباطی را بازی می کنند. به صورت شهودی توزیع یال ها در گراف های واقعی همانند گراف فیس بوک نه به

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

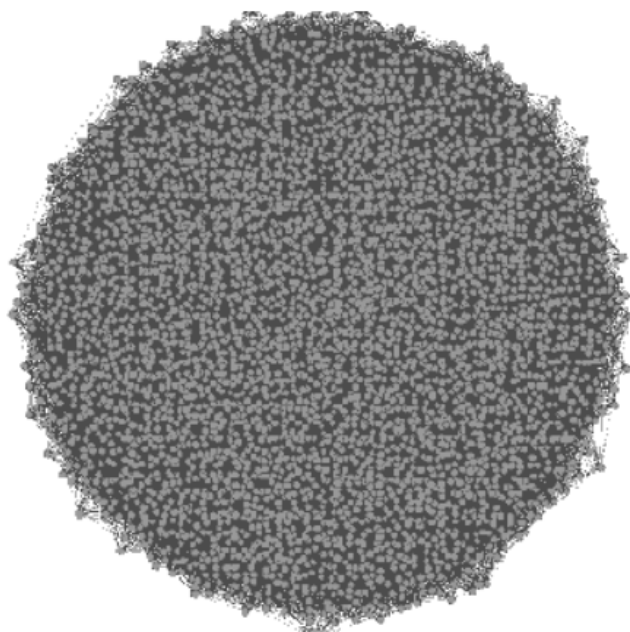
صورت سراسری که به صورتمحلی است. بدین معنی که تعداد یال های توزیع شده میان گروهی از راس ها بسیار بیشتر از تعداد یال های توزیع شده بین این گروه از راس ها با بقیه راس های گراف است. این ویژگی که در گراف های مربوط به داده های واقعی می توان مشاهده کرد، انجمن خوانده می شود. برخی از منابع انجمن، خوشه^۳ یا ماجول^۴ نیز خوانده می شود. به عبارت دیگر انجمن ها مجموعه ای از راس هاست که با احتمال بیشتری نسبت به بقیه گراف ویژگی های مشترک را به اشتراک می گذارند. انجمن ها و انجمن یابی چندین سال است که به صورت گسترده مورد مطالعه قرار می گیرد [۵۶، ۷۸]. در این پژوهش به گروه بندی که بین راس های گراف انجام می شود، تقسیم بندی، به هریک از این گروه ها در طول فرایند انجمن یابی، خوشه و پس از اتمام فرایند، انجمن می گویند. از کاربردهای انجمن یابی می توان به تبلیغات و بازاریابی اشاره کرد. از آنجایی که افراد حاضر در انجمن های تشکیل شده در یک شبکه اجتماعی به احتمال زیاد علایق مشترکی دارند، می توان با یافتن علایق آن ها از این اطلاعات به منظور تبلیغ محصولات خاص استفاده کرد. کاربردهای زیاد دیگری نیز می توان برای انجمن ها نام برد. همین کاربردهاوان در زمینه هایی همچون زیست شناسی، مهندسی کامپیوتر، اقتصاد این شاخه از علم شبکه های اجتماعی و تئوری گراف را به زمینه ای محبوب برای پژوهشگران جهت تحقیق تبدیل کرده است. برای مثال فرض کنید انجمن شامل دانش آموزان مربوط به یک مدرسه باشند این گروه از اعضای یک شبکه دارای ارتباطات بیشتری نسبت به بقیه اعضای حاضر در یک شبکه اجتماعی هستند، حال در این مدرسه دانش آموزان مربوط به رده های تحصیلی مختلف ارتباطات بیشتری دارند. برای مثال دانش آموزان مربوط به کلاس های اول ارتباطات بیشتری با هم دارند و انجمن کوچتری را تشکیل می دهند در همین انجمن کوچتر دانش آموزانی که در یک کلاس خاص (مثلاً اول یک) هستند نیز ارتباطات بیشتری با هم دارند و تشکیل انجمن کوچتری را می دهند که انجمن بزرگتر اجتماعی از این مجموعه هاست. برخی از الگوریتم های پیشنهاد شده برای انجمن یابی دقیقاً به کاوش در این ساختار سلسله مراتبی می پردازند. هرچه گام های بیشتری اجرا شوند انجمن های کوچتر را کشف می کنند به طور خلاصه می توان گفت که هدف الگوریتم های انجمن یابی کشف ساختار انجمنی، و در صورت وجود یافتن ویژگی سلسله مراتبی گراف ها تنها با اتکا به ساختار توپولوژیکی و اطلاعات ارائه شده توسط گراف است [۴]. ریشه های انجمن یابی را می توان در [۹] پی گرفت که در آن ویس^۵ و جاکوبسن^۶ برای یافتن کارگروه های یک شرکت نخستین بار این مفهوم را به کار^۱ گرفتند. اهمیت این مقاله به این خاطر است که به معرفی برخی از مهمترین الگوریتم های کشف انجمن در گراف شبکه اجتماعی می پردازد.

^۱Hub

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها



شکل ۱: زیرگرافی از گراف کامل فیسبوک با ۴۰۳۹ راس و ۸۸۲۳۴ یال



شکل ۲: گراف تصادفی ۴۰۳۹ راس و ۸۸۲۳۴ یال

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

به هر ساختار اجتماعی از افراد که بر اساس یک رابطه اجتماعی ایجاد می شود، یک شبکه اجتماعی می گوئیم. بنابراین هر شبکه اجتماعی شامل مجموعه ای از انسان ها و روابط اجتماعی بین آن هاست. لذا هر شبکه اجتماعی از دو عنصر تشکیل شده است: موجودیت های شرکت کننده در ارتباط و ارتباط بین این موجودیت ها. در علوم اجتماعی به موجودیت های شرکت کننده در ارتباط با دیگر و به ارتباطات بین این موجودیت ها رابطه گفته می شود. شبکه های اجتماعی به دو نوع برخط و برون خط تقسیم می شود. از شبکه های برون خط می توان به شبکه دوستان، شبکه همکاران و شبکه همکلاسی ها اشاره کرد. از شبکه های برخط می توان به شبکه های اجتماعی نظیر فیس بوک^۱، توییتر^۲ و گوگل پلاس^۳ اشاره کرد. شبکه های اجتماعی از قرن نوزدهم مورد توجه قرار گرفت. پژوهش ها در این حوزه از دهه چهل به بعد با تعریف ابزارهایی چون گراف اجتماعی [۲] شتاب بیشتری گرفت. در سال 1994 واسرمن^۴ با چاپ کتاب تحلیل شبکه های اجتماعی [۱] این زمینه از علم را وارد دوره جدیدی کرد، و پس از آن شبکه های اجتماعی به صورت جدی در زیرمجموعه های علوم اجتماعی و ریاضی مورد بررسی قرار گرفت. بحث های جسته و گریخته ای از سال 1960 در مورد شبکه های اجتماعی برخط به راه افتاد. نخستین شبکه های اجتماعی در سال 1997 با نام سیکس دیگرز^۵ راه اندازی شد. اما انقلاب عظیم در هزاره دوم میلادی به وقوع پیوست جایی که از سال 2002 به بعد شبکه هایی نظیر فرنداستر^۶ و اورکات^۷ و لینکدین^۸ روی وب قرار گرفتند. پدیده بزرگ دیگر شبکه اجتماعی فیس بوک بود که در سال 2004 توسط مارک زاکربرگ^۹ به دنیا معرفی شد.

۲- تجزیه و تحلیل شبکه های اجتماعی

تجزیه و تحلیل شبکه اجتماعی برای بررسی ساختار روابط اجتماعی یک گروه با هدف کشف ویژگی ها و روابط گروه یا افراد می باشد. همچنین اینطور هم می توان تعریف کرد، تجزیه و تحلیل شبکه های اجتماعی به درک رابطه بین "بازیگران" می باشد که بازیگر (گره) می تواند یک فرد، یک سازمان، یک رویداد یا یک شی باشد. امروزه تجزیه و تحلیل شبکه های اجتماعی مورد مطالعه محققان رشته هایی مانند: جامعه شناسی، ارتباطات، علوم کامپیوتر، آموزش و پرورش، اقتصاد، جرم شناسی، علم مدیریت، پزشکی، علوم سیاسی و سایر رشته ها قرار گرفته است

۳- تعریف انجمن

همانطور که در مقدمه اشاره شد، یکی از دشواری های مربوط به انجمن ها ارائه تعریفی برای آن است. تعریف انجمن ها نیازمند مقداری چشم پوشی و اختیار است، تا از دشواری های کار کاسته شود. از طرفی اصولاً تعریف دقیق انجمن ها برای کاربردهای مختلف می تواند متفاوت باشد. اگر بین هر دو راس تابعی به عنوان تابع فاصله وجود داشته باشد، می توانیم انجمن ها را به عنوان مجموعه ای از راس ها تعریف کنیم که به فاصله کمی از هم قرار دارند. این تعریف بیشتر در خوشه بندی داده ها استفاده می شود. اما همچنان که پیش تر عنوان شد انجمن ها مجموعه ای از راس هاست که تعداد یال های داخل آن مجموعه (یعنی یال هایی که هر دو راس انتهایی آن داخل انجمن باشند) از یال های آن مجموعه، با دیگر راس های گراف به مراتب بیشتر باشد. این تعریف به عنوان پایه ای برای اغلب تعاریف دیگر به حساب می آید.

¹ www.twitter.com

² www.plus.google.com

³ Wasserman

⁴ SixDegree.com

⁵ Orkut

⁶ Linkdin

⁷ MarkZuckerberg

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

۴- الگوریتم های انجمن یابی

الگوریتم های ارائه شده برای انجمن یابی را می توان به شش گروه تقسیم کرد:

1. الگوریتم های تقسیم
2. الگوریتم های براساس ماجولاریتی
3. الگوریتم های براساس شعاع طیفی مقادیر ویژه
4. الگوریتم های پویا
5. الگوریتم های براساس استنباط های آماری^۳
6. سایر روش ها^۴

هر کدام از الگوریتم های ارائه شده برای انجمن یابی در یکی از این دسته ها قرار می گیرد. در ادامه این بخش به تعریف این گروه ها خواهیم پرداخت و برای آن ها معروف ترین الگوریتم را شرح خواهیم داد

۱-۴ الگوریتم های تقسیم

فلسفه کلی این دسته از الگوریتم ها یافتن یال های بین انجمنی و حذف آن ها است. اگر تمام یال های بینانجمن ها را حذف کنیم، آن انجمن ها مولفه های جدا از هم را تشکیل می دهند معروف ترین و محبوب ترین الگوریتم که در دسته الگوریتم های تقسیم قرار می گیرد الگوریتم، گیروان نیومن [۳][۱۰] است. این الگوریتم برای یافتن یال های بین انجمنی از مفهوم مرکزیت ارتباطی استفاده می کند. ایده کلی این الگوریتم این است یال هایی که بین انجمن ها قرار می گیرند، دارای مقدار مرکزیت ارتباطی بزرگ تری هستند. این الگوریتم از چهار مرحله تشکیل می شود:

1. محاسبه مرکزیت ارتباطی هر یال،
2. یال با بزرگترین مقدار مرکزیت ارتباطی را حذف می کنیم،
3. مرکزیت ارتباطی یال ها را دوباره محاسبه می کنیم
4. به مرحله دوم باز می گردیم

نخست این الگوریتم در مقاله [۱۰] معرفی شد. با توجه به اینکه زمان محاسبه مرکزیت ارتباطی برای هر یال در هر مرحله برابر $O(n^2)$ (برای ماتریس های تنک) پیچیدگی زمانی این الگوریتم برابر $O(n^3)$ بود. در مقاله ای که نیومن^۱ و گیروان^۲ در سال 2004 [8] دو سال پس از مقاله اول منتشر کردند؛ معیاری برای یافتن بهترین تقسیمارائه دادند. در نسخه اولیه تمام نمودار سلسله مراتبی الگوریتم بدست می آمد و با برش از یک سطح انجمن هاشخص می شد اما در الگوریتم اصلاح شده، تقسیمی که بزرگترین مقدار ماجولاریتی را ارائه می داد، به عنوان تقسیم مورد نظر مشخص می کردند.

این الگوریتم با وجود بهبودهای حاصل شده پیچیدگی زمانی زیادی داشت و مهمتر از همه نمی توانست انجمن های همپوشان را کشف کند.

^۱Newman
^۲Girvan

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

۲-۴ الگوریتم های بر اساس ماجولاریتی

هرچقدر مقدار ماجولاریتی یک تقسیم به مقدار ماجولاریتی ماکزیمم نزدیک تر شود، تقسیم حاصل تقسیمخوب تری خواهد بود. همین مساله ایده اصلی این دسته از الگوریتم ها هستند. الگوریتم های مربوط به این دسته به چهار نوع زیر تقسیم می شوند

- الگوریتم تبرید شبیه سازی شده
- بهینه سازی خارجی
- بهینه سازی شعاع طیفی
- دیگر شیوه های بهینه سازی

الگوریتم های متفاوتی برای بیشینه سازی ماجولاریتی ارائه شده که به یکی از چهار دسته بالا تعلق دارند. اما از مهمترین الگوریتم های معرفی شده الگوریتم حریصانه ای است که توسط نیومن [۱۱] ارائه شد. الگوریتم با قرار دادن هر راس در یک انجمن مجزا کار خود را شروع می کند. در شروع کار هیچ یالی وجود ندارد، با افزودن یکبه یک یال ها انجمن هایی که در دو سر این یال قرار دارند در صورت افزایش ماجولاریتی تقسیم در هم ادغامی شوند. ماجولاریتی تقسیم از روی گراف کامل محاسبه می شود یعنی گرافی که یال ها به آن اضافه می شود، تنها حکم نشانگر انجمن ها را دارد. اگر افزودن یک یال ادغامی در انجمن ها به وجود نیاورد، آن یال یک یال درون انجمنی است. لذا مقدار ماجولاریتی را تغییر نخواهد داد. تعداد تقسیم های یافت شده در طول فرایند برابر n یعنی تعداد راس هاست. هرکدام از این تقسیم ها دارای یک مقدار ماجولاریتی هستند در نهایت بعد از افزودن یال ها تقسیمی که بزرگترین ماجولاریتی را دارد به عنوان خروجی مشخص می شود پیچیدگی زمانی این الگوریتم برای یک گراف تنک یعنی گراف هایی که تعداد یال های آن از تعداد کل یال های ممکن که برای n راس برابر $n(n-1)/2$ است بسیار کمتر (از مرتبه خطی) می باشد، مساوی $O(n^2)$ است.

۳-۴ الگوریتم های بر اساس شعاع طیفی مقادیر ویژه

این الگوریتم ها از شعاع طیفی مقادیر ویژه ماتریس مربوط به ماتریس های مجاورت (یعنی ماتریس هایی که بهطریقی از ماتریس مجاورت استخراج می شوند، استفاده می کنند. نخستین تحقیقات بر روی شعاع طیفی خوشه ها توسط دوناث^۱ و هوفمن^۲ [۱۲] انجام شد. در این مقاله از بردار ویژه ماتریس مجاورت برای تقسیم بندی گراف استفاده شد. در همان سال فیدلر^۳ [۱۳] نشان داد که بردار ویژه دومین مقدار ویژه کوچک ماتریس لاپلاسین به احتمال زیاد تقسیمی ارائه می دهد که مینیمم برش این تقسیم اندازه بسیار کوچک تری دارد. گراف ساده G را با n یال داریم، ماتریس لاپلاسین آن $L = (l_{i,j})_{n \times n}$ به صورت زیر تعریف می شود:

$$L = D - A$$

که در آن D ماتریس درجات راس هاست. برحسب کاربرد می توان ماتریس لاپلاسین را به صورت زیر نیز تعریف کرد

^۱Donath
^۲Hoffman
^۳Feidler
^۴Donetti
^۵Munoz

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

$$l_{i,j} = \begin{cases} \deg(v_i) & \text{اگر } i = j \\ -1 & \text{اگر } i \neq j \text{ and } v_i \text{ adjacent to } v_j \\ 0 & \text{otherwise} \end{cases}$$

تاکنون لاپلاسیان پرکاربردترین ماتریس برای الگوریتم های شعاع طیفی بوده است [۴]. دوتی^۴ و مونز^۵ [۱۴] روشی، رابریایه بردار ویژه ماتریس لاپلاسیان ارائه دادند. چون مقدار بردار ویژه مولفه ها برای راس هایی که در یک انجمن قرار دارند، مقدار نزدیکی به هم دارند، می توان با استفاده از بردارهای ویژه انجمن ها را کشف کرد. بدین معنی که اگر از m بردار ویژه استفاده کنیم، می توانیم راس ها را در یک فضای m بعدی قرار دهیم که انجمن ها به صورت راس هایی که در گروه های نزدیک به هم در این فضا قرار دارند مشخص می شوند. هرچقدر تعداد بردارهای ویژه به کار رفته بیشتر باشد، انجمن ها به صورت واضح تری در فضا مشخص می شود. الگوریتم ارائه شده توسط دوتی و مونز شامل گروه بندی نقاط و استخراج تقسیم می باشد. دوتی و مونز از خوشه بندی سلسله مراتبی، با این محدودیت که تنها خوشه هایی که حداقل یک یال بین خوشه ای در گراف اصلی دارند با هم ادغام می شوند، استفاده می کند. از بین همه تقسیم های استخراج شده، تقسیمی که بزرگ ترین مقدار مارجولاریتی را دارد به عنوان خروجی مشخص می شود. پیچیدگی زمانی این روش برابر $O(n^3)$ می باشد.

۵-۴ الگوریتم های بویا

این دسته از الگوریتم ها بصورت مستقیم بر روی گراف کار می کنند و به این طریق انجمن ها را استخراج می کنند. در اینجا به تشریح الگوریتم عابر تصادفی که توسط ژو^۱ [۱۵] ارائه شده است و در این دسته از الگوریتم ها طبقه بندی می شود، می پردازیم. عابر تصادفی که نخستین بار توسط هوگز^۲ [۱۶] معرفی شد، عابری تصادفی است که به صورت تصادفی روی گراف حرکت می کند و در هر راس با توجه به یال های موجود به هر کدام از راس های مجاور بصورت تصادفی می رود. ایده الگوریتم ژو این است که عابر تصادفی به علت چگالی زیاد یال ها در داخل انجمن زمان بیشتری را داخل انجمن مصرف خواهد کرد. ژو از عابر تصادفی برای تعریف فاصله بین دو راس استفاده کرد: فاصله d_{ij} بین دو راس u و v میانگین تعداد یال هایی است که یک عابر تصادفی برای رسیدن از u به v باید از آن ها عبور کند. راس ها نزدیک به هم احتمالاً به یک انجمن تعلق دارند و جذب کننده سراسری راس u به عنوان نزدیک ترین همسایه این راس تعریف می کند (راسی که کوچکترین مقدار d_{ij} داشته باشد). او همچنین جذب کننده محلی راس u به عنوان راسی تعریف می کند که نزدیک ترین همسایه آن باشد. دو نوع انجمن براساس جذب کننده سراسری و جذب کننده محلی وجود دارد: راس u تمام راس هایی که جذب کننده سراسری (محلی) هستند در یک انجمن قرار می گیرد و همه راس هایی که جذب کننده سراسری (محلی) آن هستند نیز در انجمنی قرار می گیرند که u در آن قرار دارد. انجمن زیرگراف کمینه است، یعنی هیچ زیرگراف کوچک تری که شرایط مورد نظر در آن صدق کند وجود ندارد. پیچیدگی زمانی این الگوریتم برابر $O(n^3)$ است.

۶-۴ الگوریتم های براساس استنباط های آماری

استنباط های آماری به منظور استخراج ویژگی های مجموعه ای از داده ها با استفاده از یک دسته مشاهدات و مقایسه آن ها با مدل های فرضی

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

تولید شده، مورد استفاده قرار می گیرند اگر مجموعه داده یک گراف باشد، مدل مجموعه ای از راس هاست ϵ که به وسیله یال هایی به هم متصل می شوند، این مدل تولید شده با توپولوژی گراف اصلی منطبق است. این دسته از الگوریتم ها از این مدل ها استفاده کرده و انجمن های موجود در گراف راپیش بینی می کند. استنباط بیزی یکی از روش های استنباط های آماری است که در مدل سازی گراف های واقعی مانند شبکه های اجتماعی به کار می رود، در این بخش به معرفی الگوریتم هاستینگ^۱ [۱۷] می پردازیم که از روش استنباط بیزی استفاده می کند. استنباط بیزی از مشاهدات به منظور تخمین احتمال درستی یک فرض استفاده می کند. استنباط بیزی شامل دو جزء است: شواهد، که عبارت است از اطلاعات D که می توان از سیستم بدست آورد؛ و یک مدل آماری با پارامتر θ استنباط بیری از محاسبه احتمال $P(D|\theta)$ که برابر احتمال مشاهده شواهد در مدل مورد نظر با پارامتر θ است، شروع می شود. هدف مشخص کردن مقدار θ ای است که مقدار $P(D|\theta)$ را حداکثر کند. در ارتباط با گراف، شواهد توسط ساختار گراف (ماتریس مجاورت یا ماتریس وزن) بدست می آید. در این مورد یک جزء دیگر علاوه بر اجزاء ذکر شده وجود دارد و آن گروه بندی گراف ها به وسیله قرار دادن راس ها داخل گروه هاست. این گروه بندی اطلاعات پنهانی است که انتظار داریم از مدل انتخابی به وسیله پارامتر بدست آوریم. در تمام روش هایی که از استنباط بیزی استفاده می کنند، هدف پیشینه کردن $P(D|\theta)$ می باشد که در آن مدل شامل ساختار گراف مشاهده شده، با مقداری محدودیت اعمال شده می باشد. در ارتباط با گراف ها، پارامتر θ توسط سه گانه $(\{q\}, \{\square\}, \{k\})$ مشخص می شود که در آن $\{q\}$ انجمن هایی که توسط انتساب راس ها بهیال ها مشخص شده را نشان می دهد θ پارامتر مدل و k تعداد خوشه هاست. الگوریتم هاستینگ از مدلی با نام مدل طراحی شده استفاده می کند که در این مدل، n راس به q گروه منتسب می شوند: راس هایی که در یک گروه قرار دارند با احتمال p_{in} هم در ارتباط هستند (یالی بینشان قرار دارد)، در حالی که راس هایی متعلق به گروه های مجزا با احتمال p_{out} هم ارتباط دارند. اگر $p_{in} > p_{out}$ گراف دارای انجمن می باشد. کلاس بندی گراف با مجموعه برچسب های $\{q_i\}$ مشخص می شود. احتمال اینکه با توجه به گراف داده شده کلاس بندی $\{q_i\}$ یک کلاس بندی مناسب، متناسب با مدل داده شده باشد برابر است با

$$p(\{q_i\}) \propto \{\exp[-\sum_{(ij)} J\delta_{q_i q_j} - \sum_{i \neq j} J'\delta_{q_i q_j}/2]\}^{-1}$$

که در آن $J = \log\{[p_{in}(1 - p_{out})]/[p_{out}(1 - p_{in})]\}$ و $J' = \log\{[(1 - p_{in})]/[(1 - p_{out})]\}$

و اولین مجموع روی نزدیک ترین همسایه ها اجرا می شود.

۴-۷ سایر روش ها

بقیه الگوریتم های ارائه شده برای یافتن انجمن ها در این دسته قرار می گیرند. از بین این الگوریتم ها، الگوریتم راقوان^۲ و همکاران [۱۸] را معرفی خواهیم کرد. راقوان و همکاران الگوریتم ساده و سریعی طراحی کردند که با نام الگوریتم گسترش برچسب ها شناخته می شود. راس ها در ابتدای کار با برچسب های منحصر به فردی از هم مجزا می شوند. در هر مرحله یک بار مورد بررسی قرار می گیرند و برچسب شان با کمک برچسب راس های اطراف مشخص می شود؛ بدین ترتیب که هر راس برچسب خود را به برچسبی که بیشترین تعداد تکرار را در بین راس های

^۱ Zhou-2Hughes

^۱ Hastings

^۲ Raghavan

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

همسایه دارد، تغییر می دهد. اگر چندین برچسب با تعداد تکرار مساوی در همسایگی راسی قرار داشته باشد، یکی به تصادف انتخاب می شود. بدینوسیله برچسب ها در گراف گسترش می یابند: اغلب برچسب ها ناپدید می شوند و برخی دیگر به برچسب غالب گراف تبدیل می شوند. در مرحله ای که دیگر تغییری رخ نمی دهد، الگوریتم پایان می یابد. به صورت ساختار پهر راس همسایه های بیشتری را در بین راس های انجمنی که در آن قرار دارد، نسبت به راس های متعلق به دیگر انجمن ها داراست. این الگوریتم احتمالا جواب های مختلفی را برای یک گراف در تکرار (الگوریتم) بدست می دهد. الگوریتم برای هر مرحله زمان $O(m)$ را مصرف می کند.

۵- نتیجه گیری

در این پژوهش تلاش کردیم الگوریتم هایی را که در بحث کاوش شبکه های اجتماعی به منظور جستجوی وجود انجمن ها بطور معمول و گسترده بکار میروند، معرفی کنیم. پس از معرفی ماجولاریتی یک دسته از الگوریتم ها معرفی شدند که به کمک آن ها می توان انجمن های یک گراف را بدست آورد. بهینه سازی مقدار ماجولاریتی یکی از پایه های اصلی این الگوریتم هاست. کاربرد دوم مقایسه مقدار ماجولاریتی بیشینه با مقدار ماجولاریتی انجمن حاصل از یک الگوریتم به منظور بررسی کارایی الگوریتم و کیفیت انجمن بدست آمده می باشد. ماجولاریتی مشکلاتی را به همراه دارد. یکی از مشکلات مربوط به بیشینه مقدار می باشد و اینکه این بیشینه مقدار تا چه اندازه می تواند مورد اعتماد باشد. اما مشکل بنیادی تر به این برمی گردد که ماجولاریتی تا چه اندازه قابلیت کشف و شناسایی انجمن های خوب را دارد این مشکل ماجولاریتی به خاطر تعریف مدل محض می باشد. نقطه ضعف مدل محض در این است که راس ها با همه راس های دیگر قادر به تعامل هستند، بدین معنا که هر قسمت از گراف همه اطلاعات مربوط به دیگر قسمت های گراف را در گستره دیدش دارد با وجود همه این موارد که به عنوان نواقص تابع کیفیت ماجولاریتی بر شمرده می شود تا کنون ماجولاریتی پر کاربردترین تابع کیفیت ارائه شده می باشد. [4]

مراجع:

- [1] Wasserman, S. (1994). Social network analysis: Methods and applications (Vol. 8). Cambridge university press.
- [2] Moreno, J. L. (1953). Who shall survive? Foundations of sociometry, group psychotherapy and socio-drama.
- [3] Newman, M. E., Girvan, M. (2004). Finding and evaluating community structure in networks. Physical review E, 69(2), 026113.
- [4] Fortunato, S. (2010). Community detection in graphs. Physics Reports, 486(3), 75-174..
- [5] Dorogovtsev, S. N., Mendes, J. F. F., Samukhin, A. N. (2000). Structure of growing networks with preferential linking. Physical Review letters, 85(21), 4633.
- [6] Krapivsky, P. L., Redner, S., Leyvraz, F. (2000). Connectivity of growing random networks. Physical review letters, 85(21), 4629.
- [7] Toivonen, R., Kovanen, L., Kivelä, M., Onnela, J. P., Saramäki, J., Kaski, K. (2009). A comparative study of social network models: Network evolution models and nodal attribute models. Social Networks, 31(4), 240-254.
- [8] Coleman, J. S. (1964). Introduction to mathematical sociology. London Free Press Glencoe.
- [9] Weiss, R. S., Jacobson, E. (1955). A method for the analysis of the structure of

دومین همایش ملی مهندسی کامپیوتر، داده های حجیم و الگوریتم ها

- complex organizations. American Sociological Review, 661-668.
- [10] Girvan, M., Newman, M. E. (2002). Community structure in social and biological networks. Proceedings of the National Academy of Sciences, 99(12), 7821-7826.
- [11] Newman, M. E. (2004). Fast algorithm for detecting community structure in networks. Physical review E, 69(6), 066133.
- [12] Donath, W. E., Hoffman, A. J. (1973). Lower bounds for the partitioning of graphs. IBM Journal of Research and Development, 17(5), 420-425.
- [13] Fiedler, M. (1973). Algebraic connectivity of graphs. Czechoslovak Mathematical Journal, 23(2), 298-305.
- [14] Donetti, L., Munoz, M. A. (2004). Detecting network communities: a new systematic and efficient algorithm. Journal of Statistical Mechanics: Theory and Experiment, (10), P10012.
- [15] Zhou, H. (2003). Distance, dissimilarity index, and network community structure. Physical review e, 67(6), 061901.
- [16] Hughes, B. D. (1996). Random walks and random environments. Oxford: Clarendon Press.
- [17] Hastings, M. B. (2006). Community detection as an inference problem. Physical Review E, 74(3), 035102.
- [18] Raghavan, U. N., Albert, R., Kumara, S. (2007). Near linear time algorithm to detect community structures in large-scale networks. Physical Review E, 76(3), 036106.