

به نام خدا

گزارش مسئله دوم
مسابقه مهارت سنجی فن آورد
بخش داده کاوی

پاییز ۹۵

مختصری از مطالب پیش رو:

- طرح مساله
- پیش پردازش
- نمودار گراف اولیه
- روش پیشنهادی
- معیار های ارزیابی
- نتایج
- گراف اجتماعات حاصل

طرح مسئله:

در این مسئله، تعداد ۲۸۶۵۷ ردیف داده از ارتباطات (تراکنش) افراد در یک شبکه اجتماعی داده شده است. در این دادگان هر ردیف داده شامل id دو فرد و وزن ارتباط آنها میباشد، میتوان برای این دادگان یک گراف وزن دار بی جهت (با توجه به غیر تکراری بودن داده ها) به اصطلاح گراف شبکه پیچیده (Complex Network) در نظر گرفت به این صورت که افراد رئوس و ارتباطات یال های گراف باشند. یا توجه به مجموع وزن های یال های گراف میتوان این گراف را به مجموعه های گسسته از راس ها خوشه بندی کرد به نحوی که مجموع این رئوس گراف اصلی را تشکیل میدهد و به این ترتیب میتوان اجتماعات کاربران این شبکه را مدل سازی کرد.

پس با یک مساله خوشه بندی گراف (Graph Clustering) از نوع درون گراف (Within-graph) و مشخصا مساله Community Structure Detection روبرو هستیم.

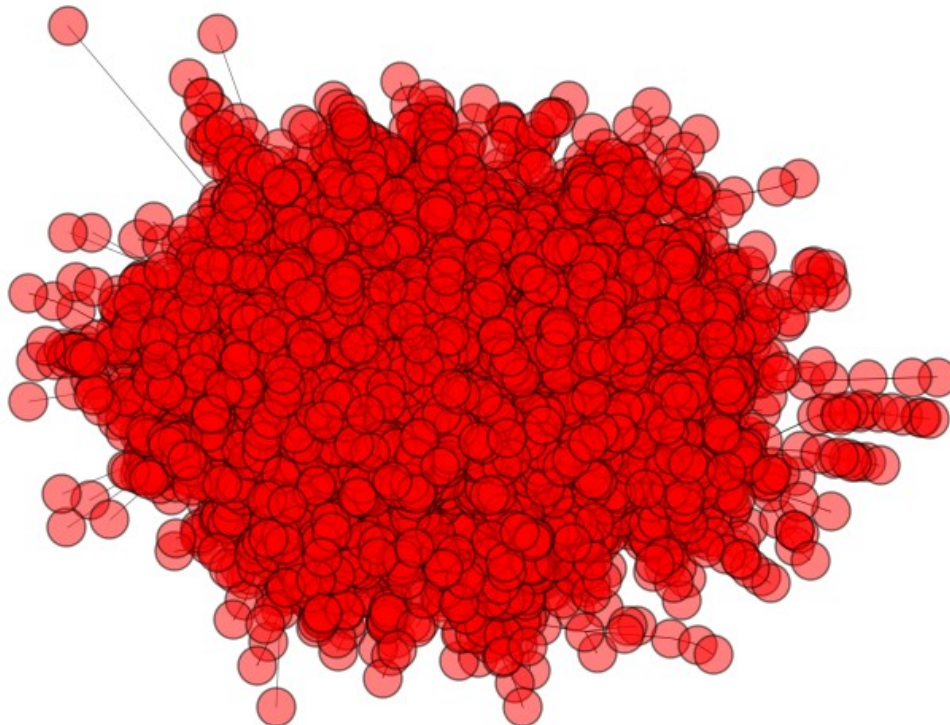
از روش های مرسوم برای حل این مساله میتوان به روش های زیر اشاره کرد:

- k-Spanning Tree
 - Shared Nearest Neighbor
 - Betweenness Centrality Based
 - Highly Connected Subgraph یا Highly Connected Components
 - Maximal Clique Enumeration
 - Kernel k-means
- در این مساله خاص از روش louvain community detection که از روش های Modularity Optimization میباشد، استفاده شده است که در ادامه توضیح داده خواهد شد.

پیش پردازش:

در این مساله از آنجایی که تمام داده ها **Numeric** یا عددی می باشند، داده ها نیاز به پیش پردازش زیادی ندارند. پس از تبدیل فایل دادگان به فرمت مناسب استفاده، ستون های مبدا (**Source**)، مقصد (**Destination**) و وزن (**weight**) به ترتیب برای تشکیل گراف استفاده شده است. این کار توسط کتابخانه **NetworkX** و به زبان **Python** انجام شده است.

نمودار گراف اولیه:



روش پیشنهادی:

در روش های **Modularity Optimization** برای افزایش وزن دار به اجتماعی از رؤس، یک تابع هزینه مناسب تعریف میشود که بر اساس آن میتوان گراف را خوشه بندی کرد. برای این کار گراف با ماکزیمم ماژولاریتی توسط متد **best_partition** از کتابخانه **Community** محاسبه میشود. گراف حاصل بیشترین ماژولاریتی را مطابق الگوریتم **Louvain** دارد.

معیار های ارزیابی:

معیارهای گوناگونی برای ارزیابی خوشه بندی مناسب یک گراف وزن دار وجود دارد که در ادامه، به معیار انتخاب شده و دلیل انتخاب آن می پردازیم.

دو معیار مهم که برای ارزیابی مناسب بودن یک خوشه استفاده می شوند، **چگالی درون-خوشه ای و بین-خوشه ای** هستند. در یک خوشه بندی مناسب سعی بر این است که چگالی درون-خوشه ای بیشینه شود و چگالی بین-خوشه ای کمینه. این دو معیار طبق فرمول زیر در یک گراف قابل محاسبه هستند:

$$\delta_{\text{int}}(C) = \frac{\# \text{internal edges of } C}{n_c(n_c - 1) / 2}$$

$$\delta_{\text{ext}}(C) = \frac{\# \text{inter-cluster edges of } C}{n_c(n - n_c)}$$

• چگالی درون-خوشه ای :

نشان دهنده نسبت تعداد یال های در یک خوشه به تعداد کل یال های ممکن در آن خوشه است. همچنین چگالی بین-خوشه ای نشان دهنده تعداد یال های بین یک خوشه با سایر رؤس گراف، نسبت به کل تعداد یال های ممکن در گراف است. یک معیار مناسب که می توان تابع هزینه را بر اساس آن تعریف کرد تفاضل میان فاصله ی درون-خوشه ای و بین-خوشه ای است.

دو معیار بالا همان طور که مشخص است برای گراف های ساده و غیروزن دار کاربرد دارند. در صورت وزن دار بودن گراف این معیارها، وزن ها را در نظر نمی گیرند و تنها وجود یا عدم وجود یال بین رأس ها را مورد بررسی قرار می دهند. بنابراین نیاز به تغییری در معیار بالا برای قابل استفاده بودن در مسئله خوشه بندی گراف های وزن دار است.

معیاری که برای یک خوشه‌بندی در یک گراف وزن‌دار می‌تواند مناسب باشد آن است که مجموع اوزان یال‌ها درون یک خوشه یا اجتماع بیشینه و مجموع اوزان یک خوشه یا اجتماع با خوشه‌ها و اجتماعات دیگر کمینه شود. بنابراین معیار شهودی، معیار **ماژولاریتی**^۱ که در بخش بعد به آن می‌پردازیم به عنوان معیار برای این مسئله انتخاب شده است.

• معیار خوشه‌بندی ماژولاریتی:

یک خوشه‌بندی مناسب در یک گراف، لزوماً خوشه‌بندی نیست که در آن تعداد یال‌ها بین خوشه‌ها کمترین مقدار ممکن باشد، چرا که اگر به طور مثال یک رأس تنها در یک دسته باشد و تمام رئوس دیگر در دسته‌ای دیگر، آن‌گاه بر طبق این معیار یک جواب بهینه داریم که آشکارا نامناسب است. یک خوشه‌بندی مناسب آن خوشه‌بندی است که **تعداد یال‌ها میان دو خوشه‌ی مختلف، کمتر از مقدار مورد انتظار باشد**. معیار ماژولاریتی نیز بر این اساس تعریف می‌شود. در واقع ماژولاریتی در یک خوشه از تفریق تعداد یال‌های میان رئوس از تعداد یال‌های مورد انتظار در یک حالت تصادفی حاصل می‌شود. روش دقیق به دست آوردن ماژولاریتی در زیر توضیح داده شده است.

با در نظر گرفتن ماتریس مجاورت گراف مورد نظر، منظور از A_{ij} وزن یال بین دو رأس i و j است. همچنین در صورتی که حالت تصادفی را در نظر بگیریم، همچنین K_i را درجه‌ی رأس K در نظر بگیرید. m نیز تعداد کل یال‌های گراف است. بنابراین ماژولاریتی برای یک خوشه می‌تواند به صورت زیر تعریف شود:

$$A_{vw} - \frac{k_v k_w}{2m}$$

با جمع‌بستن تمامی این مقادیر به ازای هر دو رأس، به معادله‌ی زیر می‌رسیم:

$$Q = \frac{1}{2m} \sum_{uv} [A_{uv} - \frac{k_v k_u}{2m}] \delta(c_v, c_u)$$

مقدار تابع دلتا یک است اگر دو رأس در یک خوشه یا اجتماع باشند وگرنه صفر است. حال هدف ما کسب‌موم

کردن مقدار Q است. مقدار Q عددیست در بازه‌ی $[1, 0]$ و هرچه قدر به ۱ نزدیک‌تر باشد، نشان‌دهنده‌ی دسته‌بندی بهتر است.

بنابراین مسئله‌ی اولیه تبدیل به بیشینه‌کردن ماژولاریتی می‌شود. برای این کار الگوریتم‌های گوناگونی وجود دارد که در این مسئله از الگوریتم **Louvain** استفاده شده است.

نتایج:

تعداد اجتماعات به دست آمده به طوری که ماژولاریتی حداکثر شود، با استفاده از الگوریتم **Louvain** برابر ۱۶۸ خوشه بود. در فایل **output.csv** تعلق هر یک از منابع به هر کدام از خوشه‌ها به عنوان پاسخ نهایی آمده است. این فایل دارای دو ستون است که یکی شماره‌ی منبع و دیگری شماره‌ی خوشه‌ای است که آن منبع به آن تعلق دارد.

برای ارزیابی پاسخ، دو روش انتخاب شده‌اند. یکی محاسبه‌ی میزان ماژولاریتی برای تقسیم‌بندی مذکور است که مقدار آن برابر با **Modularity = 0.971198961478** است و نشان‌دهنده‌ی آن است که اجتماعات یافته شده، دارای بیشترین ارتباط درون خود، و کمترین ارتباط ممکن نسبت به حالت مورد انتظار هستند. این عدد به این جهت که به یک نزدیک است، نشان‌دهنده‌ی خوشه‌بندی مناسب بر اساس معیار ماژولاریتی است.

همچنین گراف القایی دیگری نیز روی اجتماعات به دست آمده تعریف شده است. رأس‌های این گراف القایی، اجتماعات به دست آمده هستند و یال‌های بین هر دو رأس نشان‌دهنده‌ی مجموع اوزانیست که بین آن دو اجتماع در گراف اصلی موجود است. بدیهی است که برای هر سه تایی (i, j, k) اگر $i = j$ باشد، میزان k که نشان‌دهنده‌ی مجموع اوزان بین دو اجتماع است، باید مقدار بزرگی باشد. چرا که در صورت برابری i و j در واقع اوزان بین منابع یک اجتماع مدنظر است که باید عدد بزرگی باشد. در صورتی که i و j برابر نباشند مقدار k باید تا حد امکان کوچک باشد، چرا که مجموع اوزان میان دو خوشه‌ی متمایز هر چه قدر کمتر باشد، خوشه‌بندی بهتر انجام گرفته است. مجموع اوزان میان خوشه‌ها در فایل **induced_graph.csv** قابل مشاهده است. به طور مثال مجموع اوزان بین منابع در میان اعضای خوشه‌ی اول، مطابق اطلاعات فایل، برابر ۴۰۳۰ است، و مجموع اوزان میان منابع خوشه‌ی اول با منابع خوشه‌ی شصتم، مطابق فایل برابر با ۱۰ است.

یکی دیگر از مزایای الگوریتم **Louvain** که بر اساس معیار ماژولاریتی عمل می‌کند آن است که در هر خوشه، هر تعداد از رئوس می‌توانند قرار بگیرند و اندازه‌ی هر کدام از خوشه‌ها (تعداد رئوس در یک اجتماع) لزوماً نباید با یکدیگر برابر باشند. بنابراین خوشه‌های گوناگون اندازه‌های گوناگونی دارند که در فایل **number_of_nodes.csv** اندازه‌ی هر کدام از خوشه‌ها قابل مشاهده است.

گراف اجتماعات حاصل:

