

بررسی اجتماعات وب و الگوریتم‌های استخراج آن‌ها

توسط

مهیار سالک

رساله ارائه شده به عنوان بخشی از ملزومات برای دریافت درجه
کارشناسی نرم‌افزار کامپیوتر

زیر نظر

دکتر محمد قدسی

مرداد ماه ۱۳۸۴

دانشکده مهندسی کامپیوتر

دانشگاه صنعتی شریف

تهران

تقدیم به

پدر و مادرم

قدردانی

ازاستادم، آقای دکتر محمد قدسی که مرا به علوم نظری کامپیوتر علاقمند و همواره مرا از راهنمایی های خود بهره مند ساختند، تشکر فراوان می کنم و شایستگی ایشان را در تمامی وجوه زندگانی می ستایم. به دوست و همراهم مرتضی منعمی که مرا به این مبحث معرفی و تشویق کردند امتنان می ورزم. از آقای دکتر David Kempe که مرا با مسائل شبکه های اجتماعی و سایر مسائل مرتبط آشنا کردند و مرا در انجام این پروژه راهنمایی نمودند، بسیار سپاسگزارم. بدون کمک دوستم، محمد درخشانی، نوشتار و ویرایش این پایان نامه به زمان و هزینه بیشتری نیاز داشت. کمکش را قدر می دانم. و در خاتمه به خانواده و دوستانم به خاطر با من بودن ها ممنونم.

بررسی اجتماعات وب و الگوریتم‌های استخراج آن‌ها

چکیده

منظور از اجتماعات آن دسته از موجودیت‌های شبکه وب است که ارتباط قوی با یکدیگر دارند. این موجودیت‌ها بصورت صفحات و میزبان‌ها خود را نشان می‌دهند. برای تشخیص اجتماعات وب باید روی ساختار پیوندی بین موجودیت‌ها در گراف القایی وب مطالعه کنیم. شناخت اجتماعات وب به ما کمک می‌کند که از طریق ساختار پیوندی وب اطلاعات قابل اعتمادی در مورد محتویات صفحات وب بدست آوریم.

با داشتن گراف $G(V, E)$ که V نمایشگر صفحات وب و E نمایانگر پیوندهای بین آنهاست مساله، یافتن زیرگراف‌هایی از G است که یا بصورت هاب و مرجعی باشند (دیدگاه HITS) یا تراکم پیوندی بالا داشته باشند.

در این پایان‌نامه خصوصیات وب و مدل کردن آن بصورت یک گراف را شرح می‌دهیم. به بیان روش‌های سنتی و قدیمی‌تر یافتن اجتماعات می‌پردازیم. روش درخت برش کمینه را بررسی می‌کنیم. در ادامه روش Charikar را شرح داده و آن را روی یک مجموعه داده از صفحات وب می‌آزماییم. در انتها به جمع بندی می‌پردازیم.

واژه‌های کلیدی: اجتماع، هاب، مرجع، شاره، برش

فهرست

۷	۱ مقدمه
۷	۱-۱ اجتماعات داخل وب چه معنا و مفهومی دارند؟
۷	۱-۱-۱ مسائل مشابه
۸	۲-۱-۱ کاربردها
۹	۲-۱ اجتماعات وب به چه شکلی ظاهر می شوند؟
۹	۱-۲-۱ نگاه هاب و مرجع
۱۱	۲-۲-۱ دیدن اجتماعات بصورت یک بخش متراکم یالی در گراف وب
۱۲	۲ ساختار گرافی وب در مقیاس صفحات و میزبان ها
۱۲	۱-۲ مقدمه
۱۳	۲-۲ جستجوی عنوان
۱۴	۱-۲-۲ پارادایم حذف و تولید

۲ فهرست

۱۵ چند اندازه گیری ۳-۲

۱۵ توزیع درجه ۱-۳-۲

۱۸ اندازه قطر وب ۴-۲

۱۹ تحلیل صفحات میزبان ۵-۲

۳ الگوریتم HITS و شماره در استخراج اجتماعات وب

۲۱ الگوریتم HITS و هسته های دوبخشی ۱-۳

۲۴ شباهت HITS و Pagerank ۱-۱-۳

۲۴ سازگارشدن الگوریتم Pagerank برای تشخیص اجتماعات ۲-۱-۳

۲۵ روش شماره بیشینه در بدست آوردن اجتماعات وب ۲-۳

۲۵ مساله شماره بیشینه ۱-۲-۳

۲۷ تعریف اجتماعات از دیدگاه شماره ای و الگوریتم تشخیص آن ۲-۲-۳

۲۹ تحلیل ۳-۲-۳

۳۰ بحث و مقایسه ۳-۳

۴ الگوریتم های دسته بندی با درخت های برش کمینه

۳۲ درخت های برش کمینه ۱-۴

۳۳ الگوریتم ساختن درخت ۱-۱-۴

۳۷ الگوریتم های دسته بندی و استخراج اجتماعات ۲-۴

۳۷	الگوریتم	۱-۲-۴
۳۹	دسته‌بندی سلسله‌مراتبی	۲-۲-۴
۳۹	نرمالسازی یالی	۳-۲-۴
۴۰	کاربرد در WWW	۴-۲-۴
۴۱	الگوریتم دسته‌بندی اجتماعات	۵-۲-۴
۴۱	اجتماع ازوتریک	۶-۲-۴
۴۲	تعمیم تعریف اجتماع Esoteric	۷-۲-۴
۴۳	جستجو برای اجتماعات	۳-۴
۴۵	تسهیل شرایط	۱-۳-۴

۵ الگوریتم‌های تقریبی برای پیدا کردن اجتماعات

۴۷	گراف بدون جهت	۱-۵
۴۸	محاسبه $f(G)$	۱-۱-۵
۵۰	الگوریتم تقریبی با ضریب ۲ برای محاسبه $f(G)$	۲-۵
۵۲	پیاده‌سازی الگوریتم	۳-۵
۵۴	انتخاب d^{\max} بعنوان حد بالایی	۴-۵
۵۴	مساله‌ی معادل برای گراف جهت دار	۵-۵
۵۶	الگوریتم حریرانه برای محاسبه $d(G)$	۶-۵

۶ نتایج تجربی روی الگوریتم تقریبی Charikar

- ۵۸ الگوریتم تقریبی استخراج اجتماعات وب ۱-۶
- ۵۸ اگر کل گراف درحافظه اصلی برنامه گنجانیده نشود ۱-۱-۶
- ۵۸ اگر بخواهیم همه اجتماعات را شناسایی کنیم ۲-۱-۶
- ۵۹ داده های واقعی ۲-۶
- ۶۱ تحلیل نتایج ۳-۶

۷ نتیجه گیری

۶۳

۸ واژه نامه

۶۵

فهرست شکل‌ها

۱۰	هاب و مرجع	۱-۱
۱۵	توزیع درجات ورودی	۱-۲
۱۶	توزیع درجات خروجی	۲-۲
۱۷	نمایی از وب	۳-۲
۲۲	$K_{۲۴}$ گراف دوبخشی	۱-۳
۲۶	یک شبکه شماره	۲-۳
۳۱	چند گراف مشابه	۳-۳
۳۴	الگوریتم ساختن درخت برش کمینه	۱-۴

۳۸	تشخیص کانون در درخت	۲-۴
۳۹	دسته بندی سلسله مراتبی	۳-۴
۴۳	دو دسته بندی غیر قابل کاهش متفاوت	۴-۴
۵۳	گراف ناهنجار	۱-۵
۶۰	خروجی الگوریتم تقریبی	۱-۶
۶۱	خروجی الگوریتم تقریبی در دور دوم	۲-۶

فصل ۱

مقدمه

این پایان نامه به بررسی روش های تشخیص اجتماعات موجود در شبکه جهانی وب می پردازد. منظور از اجتماعات آن دسته از موجودیت های شبکه وب است که ارتباط قوی با یکدیگر دارند. این موجودیت ها بصورت صفحات و میزبان ها خود را نشان می دهند. اساس وجود ارتباط بین این موجودیت ها بصورت ارجاع دادن یکی دیگری را بیان می شود. اصولاً وجود پیوند بین دو صفحه نشان دهنده ارتباط معنایی بین آن دو صفحه می باشد. می توان انتظار داشت که صفحه ارجاع دهنده به یک موضوع معنایی در صفحه ارجاع شونده علاقمند است. پس تا این لحظه مشخص است که برای تشخیص اجتماعات وب باید روی ساختار پیوندی بین موجودیت ها مطالعه کنیم.

۱-۱ اجتماعات داخل وب چه معنا و مفهومی دارند؟

۱-۱-۱ مسائل مشابه

مساله تشخیص اجتماعات یک مساله جدید نیست. اصولاً مطالعات مشابهی روی سایر شبکه ها انجام شده است. بعنوان مثال شبکه ارجاعات کتب از جمله موضوع مطالعات گذشتگان بوده است. یک فرق اساسی که بین این شبکه با شبکه وب می باشد ماهیت زنده و متغیر سایت های وب در

مقابل موجودیت ایستا و همیشگی کتاب است. بدین لحاظ است که مساله زمان در مورد سایت های وب محلّیت ندارد. چه بسا صفحات پیشکسوت و قدیمی تری در وب یافت می شوند که به یک تازه وارد رجوع می دهند.

شبکه دومی که شاید بیشتر شبیه مساله ما باشد شبکه های اجتماعی^۱ می باشند. انواع مختلفی از اینگونه شبکه ها که هر کدام نوع خاصی (اجتماعی، اقتصادی و سیاسی ...) از ارتباط بین اشخاص و موجودیت های انسانی را بیان می کنند وجود دارد. اینگونه شبکه ها که اخیراً موضوع بسیار جذابی برای تحقیقات شده اند برای تحلیل و ارائه راه حل های انسانی مورد استفاده قرار می گیرند. بعنوان مثال در تحلیل نژادها و فرهنگ های قبایل بومی افریقا پیدا کردن موجودیت هایی که به عنوان پل در بین دو قوم با سیستم های بسته ایفای نقش می کنند مساله جالبی است. البته مطالعه اینگونه شبکه ها لازمه تحلیل های گسترده ای در زمینه های انسانی حاکم بر آنها علاوه بر مدل های ریاضی پیشنهادی برای آنهاست.

۲-۱-۱ کاربردها

به بحث وب باز گردیم. بازیابی اجتماعات وب کاربردهای فراوانی می تواند داشته باشد. بطور کلی تحلیل معنایی اطلاعات موجودیت های وب کاری سخت و زمان بر است و احتیاج به یک قاضی انسانی دارد. اجتماعات وب از آن جهت می توانند به ما کمک کنند چرا که بهترین منبع برای شناختن یک چیز آن چیزهایی هستند که خود به او ارجاع داده اند یا از او ارجاع گرفته اند. پس شناخت اجتماعات وب به ما کمک می کند که از طریق ساختار پیوندی وب اطلاعات قابل اعتمادی در مورد محتویات صفحات وب بدست آوریم.

چنین اطلاعاتی می تواند در بهبود کارایی موتورهای جستجو بسیار موثر واقع شود. کاربردهای متفاوت دیگری نیز می توان براساس حیطه برای چنین اطلاعاتی قائل شد. بعنوان مثال بسیاری از کاربردهای سیستم های توزیع شده وب و عملیات I/O و ... با استفاده از شناسایی اجزاء بسیار به هم مرتبط وب یا همان اجتماعات بسیار بهبود می یابد.

^۱ Social networks

۲-۱ اجتماعات وب به چه شکلی ظاهر می شوند؟

تا اینجا می دانیم اگر بخواهیم بدنبال یک اجتماع در وب بگردیم کاندیداهای ما آن دسته از موجودیت ها هستند که پیوندهای بسیاری به هم دارند. همین دید ساده سرچشمه تعاریف ریاضی متفاوتی است که برای این موجودیت معنایی داده شده اند.

اصولا دو فاکتور برای ارائه یک مدل ریاضی برای اجتماعات وب در نظر گرفته می شود.

(۱) تطابق با دید مبتنی بر وجود ارتباط قوی و ارجاع کننده و ارجاع شوندگی بین موجودیت ها

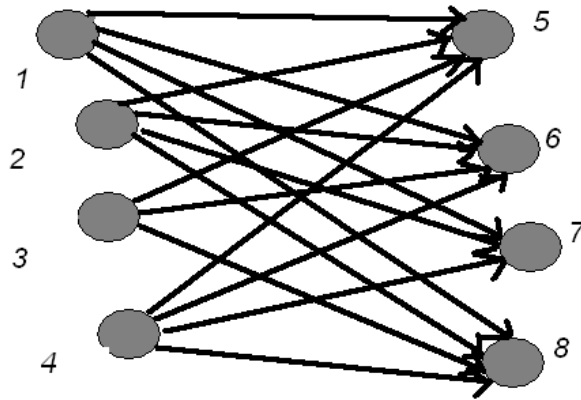
(۲) سادگی برای تشخیص از طریق راه حل مدنظر

هم این پایان نامه بررسی و مقایسه راه های مختلف تشخیص این اجتماعات است. در این پایان نامه سعی شده است بیشتر به بررسی تفصیلی راه های جدیدتر و شاید بالقوه در تشخیص اجتماعات پردازد. الگوریتم HITS و روش شبکه شاره Flake بطور اجمالی مورد بررسی قرار می گیرد و بیشتر به بیان راه حل درخت برش کمینه و الگوریتم تقریبی Charikar در شناسایی قسمت های شلوغ یک گراف و کاربرد آن در پیدا کردن اجتماعات پردازیم. به علاوه اینکه روش الگوریتم تقریبی را روی داده های حقیقی می سنجیم.

بطور کلی در روش هایی که برای شناسایی اجتماعات مطرح شده است همواره این واقعیت را مد نظر داشته اند که هیچ لزومی برای بیان یک تعریف دقیق از اجتماعات و دادن بها برای تخطی نکردن از چنین تعریفی وجود ندارد. در اکثر روش ها یا تعریف بسیار دقیقی از این موجودیت نشده است و یا اگر شده همواره یک شعاع حول آن تعریف دقیق بعنوان مجموعه جواب برگردانیده شده است. در اینجا برداشت های مختلفی که از مفهوم اجتماع در هر کدام از روش ها شده است را می آوریم.

۱-۲-۱ نگاه هاب و مرجع

اگر مفهوم ارجاع دادن را در نظر بگیریم ارجاع دهنده و ارجاع شونده دو جزء اصلی می باشند. معمولا ارجاع دهنده به یک مطلب در سایت ارجاع شونده علاقمند است یا استناد می کند. بدین



هاب و مرجع 1-1

ترتیب مجموعه‌ای از سایت‌های محبوب و مرجع یا authority و مجموعه‌ای از سایت‌های هوادار یا هاب تشکیل می‌شود. شکل ۱-۱ به بیان این مطلب می‌پردازد.

شکل ۱-۱ ساختار یک گراف دوبخشی را یادآوری می‌کند. بنابراین در اجتماعات HITS مابدنبال چنین زیرگراف‌هایی در گراف وب می‌گردیم.

پس بصورت دقیق اگر $G(V, E)$ گراف وب با V مجموعه گره‌ها که نمایانگر صفحات وب و E مجموعه یال‌ها نمایانگر ارجاعات یا پیوندهای بین صفحات هستند یک هسته دو بخشی در وب یک زیرگراف است که شامل یک گراف دو بخشی k در l می‌شود. همانطور که گفته شد به صفحات ارجاع دهنده هاب و به صفحات ارجاع شونده در هسته دو بخشی، مرجع^۲ گفته می‌شود. اجتماعاتی که HITS پیدا می‌کند به اجتماعات دو بخشی^۳ مشهور هستند.

^۲ authority

^۳ Bipartite Communities

۲-۲-۱ دیدن اجتماعات بصورت یک بخش متراکم یالی در گراف وب

اینگونه برداشت از اجتماعات که به اجتماعات متراکم یا شلوغ معروف است توسط الگوریتم های شماره و سایر الگوریتم های بررسی شده در این پایان نامه مدنظر گرفته شده است. در اینگونه اجتماعات انتظار عمومی وجود لینک های بسیار بین اعضای یک اجتماع نسبت به سایر موجودیت های وب است.

بطور دقیق تر یک اجتماع از نظر Flake یک زیرمجموعه راسی X از مجموعه V است بطوری که برای هر راس v از X تعداد یال های v به رئوس داخل X حداقل بتعداد یالهای خارج از X باشد. بیان دیگر این تعریف این است که حداقل ۵۰ درصد یال های v داخل X باشند. با این بیان دوم به راحتی می توان یک تعمیم برای این تعریف در نظر گرفت و آن پارامتر کردن عدد ثابت ۵۰ درصد است. البته Flake و بعدا Tsioutsoulouklis این مساله را از طریق شبکه شماره و با همان ثابت ۵۰ درصد حل کرده اند. پس یک پیشنهاد خوب تحقیق در مورد الگوریتم حل این مساله بصورت پارامتریزه می باشد.

دیدنی که از کار Charikar از اجتماعات بدست می آوریم هم منطبق بر اجتماعات متراکم است. در حقیقت ما به دنبال آن زیر گراف از گراف وب به عنوان یک اجتماع هستیم که بیشترین نسبت یال به راس را دارد. در آینده از این نسبت به عنوان چگالی گراف G یاد می کنیم. روش دیگری نیز بعنوان PageRank از طریق استفاده از یک قدم زن تصادفی پیشنهاد شده است. انتظار این است که یک قدم زن تصادفی خوب بیشتر روی رئوس یک اجتماع قدم زند تا رئوس خارج؛ بخصوص هنگامی که شروع قدم زدن از یک راس داخل یک اجتماع باشد.

این پایان نامه در ادامه بدین صورت منظم شده است. در ابتدا خصوصیات وب و مدل کردن آن بصورت یک گراف را شرح می دهیم. در فصل ۳ به بیان روش های سنتی و قدیمی تر یافتن اجتماعات می پردازیم. در فصل ۴ روش درخت برش کمینه را بررسی می کنیم در فصل ۵ روش Charikar را شرح داده و آن را می آزماییم. در انتها به جمع بندی می پردازیم.

ساختار گرافی وب در مقیاس صفحات و میزبان ها

۱-۲ مقدمه

نگاه کردن به وب به عنوان یک گراف جهت دار مارا قادر به بسیاری از تحلیل های معنایی می کند که بدون استفاده از چنین دیدی باید بهای فراوانی برای چنین اطلاعاتی پردازیم . این نگاه می تواند در دو واحد صفحات و میزبان ها صورت گیرد . الگوریتم های متفاوتی مانند HTS مبتنی بر چنین دید گرافی می باشد . در این فصل به بررسی ساختار غول آسای وب می پردازیم و روش های مدل کردن وب به عنوان یک گراف را می بینیم .

دلایل بسیاری وجود دارد که ما علاقمندیم چنین تحلیلی را روی وب انجام دهیم از جمله دلایل اقتصادی، ریاضیاتی، اجتماعی و .. بطور کلی کارکردن با وب بعنوان یک داده ساختار پایگاه داده ای و نادیده گرفتن اطلاعات موثری که در پیوندهای بین صفحات وجود دارد کارعقلانه ای نمی باشد . استفاده از ساختار گرافی وب به ما کمک می کند اطلاعات فراوانی در مورد محتوای صفحات و میزبان های وب پیدا کنیم . بسیاری از روش های جستجوی اجتماعات وب از چنین ساختاری بهره می برند . ساختار گراف وب بهره برداری به منظور بهبود جستجو در موتورهای

جستجو شده است. [۱، ۳]

بطور کلی اگر موضوع صفحه v برای ما ناشناخته و یا غیرقابل دسترس باشد با استفاده از مجموعه صفحاتی که v به آنها اشاره کرده و یا مجموعه صفحاتی که به v اشاره می کند می توانیم اظهارات دقیقی در مورد v و محتوای آن بکنیم. توزیع توانی به نظر می رسد توزیع مناسبی برای پیوند بین صفحات باشد. جالب است که توزیع های مشابهی برای سایر موارد مشابه مانند مرجع دادن مقالات علمی توسط [۴] Gilbert پیدا شده است. Mendelzon در [۷] به بحث درباره این موضوع پرداختند که واسط تقاضاهای SQL به پایگاه داده ها مناسب برای کشف خصوصیات جالب نهفته در ساختار گرافی وب نمی باشد. او و Wood زبان $G+$ را پیشنهاد دادند به عنوان زبانی با قابلیت بیان بالاتر نسبت به SQL برای تقاضاهای متناسب با ساختار گرافی وب کارهای مشابهی در نگاه کردن وب بصورت یک پایگاه داده ی نیمه ساختار یافته انجام شده که در نهایت به برداشت مختلط از وب بصورت رابطه ای و گرافی منجر شد.

۲-۲ جستجوی عنوان

مساله جستجوی عنوان بشکل زیر تعریف می شود: اگر یک عنوان برای جستجو داده شود (که این عنوان بصورت یک Query مطرح می شود) بعنوان خروجی صفحاتی را برگردان که در مورد آن عنوان دارای کیفیت بالایی باشند.

چنین مساله ای کاربرد هر روزه ی موتورهای جستجو است. هر کاربری که حتی برای اندک باری با موتورهای جستجو کار کرده باشد می داند که این مساله در موتورهای جستجو بصورت بهینه حل نشده است. بطور کلی برای هر عنوان خاص مجموعه ای از صفحات مرجع اطلاعاتی (authority) و صفحات بسیاری که به آنها اشاره دارند که ما از آنها به هاب نام می بریم وجود دارند

الگوریتم HITS در حقیقت به جستجوی چنین ساختارهایی می پردازد. بطور خلاصه ما به هر صفحه دو وزن نامنفی (x_p, y_p) اختصاص می دهیم که x_p وزن مرجعیت و y_p وزن اشاره دهی است. ایده اصلی این است که اگر صفحه ای با مقدار بسیاری اشاره ده با وزن بالا اشاره شد ما تمایل به

افزایش اعتبار مرجعیت آن داریم و بالعکس.

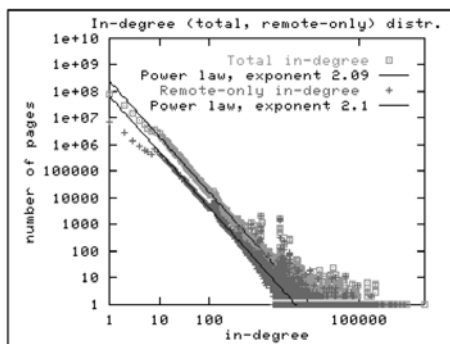
$$x_p = \sum_{q|q \rightarrow p} y_q$$

$$y_p = \sum_{q|p \rightarrow q} x_q$$

در حقیقت این درک ما را منجر به یافتن گراف‌های دارای زیر گراف‌های دوبخشی در وب می‌کند. پس الگوریتم ما بدین فرض کار می‌کند که برای یک عنوان خوش بیان در وب یک هسته در گراف وب وجود دارد. به جستجو و شمارش چنین هسته‌های دو در وب عمل *trawling* گفته می‌شود. نتایج تجربی نشان‌دهنده این واقعیت است که چنین الگوریتم لازم است که دسترسی تصادفی به یال‌های داخل گراف داشته باشد که این مستلزم نگاه دارای کل گراف در داخل حافظه اصلی است که مانع از سرباره سنگین خواندن از حافظه جانبی برای هر بار دسترسی شود. برای پیدا کردن چنین هسته‌هایی پارادایم حذف و تولید توسط [۹] پیشنهاد شده است.

۲-۲-۱ پارادایم حذف و تولید

در قسمت حذف بعنوان مثال اگر ما دنبال هسته $C_{۴,۴}$ هستیم کلیه گره‌های با درجه ورودی ۳ یا کمتر نمی‌توانند در سمت راست چنین هسته‌ای باشند پس کلیه یال‌هایی که به چنین گره‌هایی ختم می‌شوند قابل حذف می‌باشند با حذف چنین یال‌هایی مکرراً می‌توانیم عمل حذف را انجام دهیم. عمل تولید نقطه مقابل حذف می‌باشد فرض کنید در پیدا کردن $C_{۴,۴}$ گره u درجه ورودی ۴ را داشته باشد در این صورت این گره تنها در صورتی متعلق به یک هسته بصورت $C_{۴,۴}$ می‌باشد اگر و تنها اگر هر ۴ گره که به u اشاره دارد حداقل ۴ همسایه داشته باشد پس چنین فیلتری حالت صلب و اثبات دارد. اینگونه اعمال بصورت خطی با سایز گراف انجام می‌شوند.



توزیع درجات ورودی (شکل از [۶]) : 2-1

۳-۲ چند اندازه گیری

در این قسمت ما به بررسی اندازه گیری های انجام شده روی مشخصه های گراف وب می پردازیم. نتایج توزیع درجه در قسمت اول از کارهای Kumar و سایرین در [۸] Albert و سایرین در [۲] بدست آمده است. در نهایت تحلیل جزء همبند و تحلیل قطری از کارهای Broder بدست آمده است.

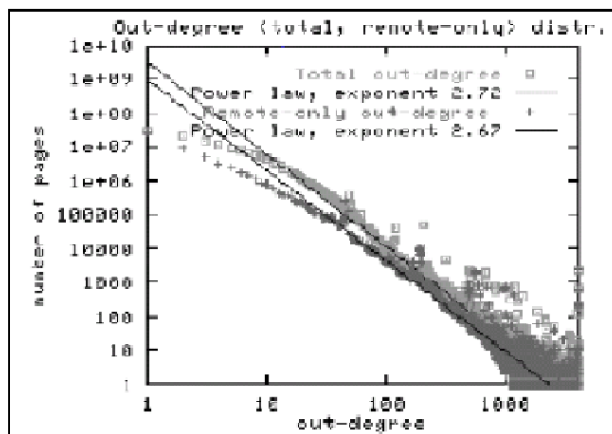
۱-۳-۲ توزیع درجه

با بررسی درجه ورودی و خروجی گره های گراف وب شروع می کنیم. کارهای اولیه kumar در [۸] به مشاهده وجود توزیع توانی در درجات ورودی بود: کسری از صفحات وب یا درجه ورودی i متناسب با $1/i^x$ برای عدد $x > 1$ می باشد. در ادامه کارهای Albert در [۲] چنین نظری را تایید کرد و یک مقدار مناسب برای x را برابر ۲.۱ تخمین زد.

نمودار log-log (با محور منفی شده x ها) در ۲-۱ از توزیع درجات ورودی می باشد^۱. همانطور که

مشاهده می شود مقدار ۲.۱ به بهترین شکل به داده های ورودی ما می خورد.

[۶]^۱



توزیع درجات خروجی (شکل از [۶]): 2-2

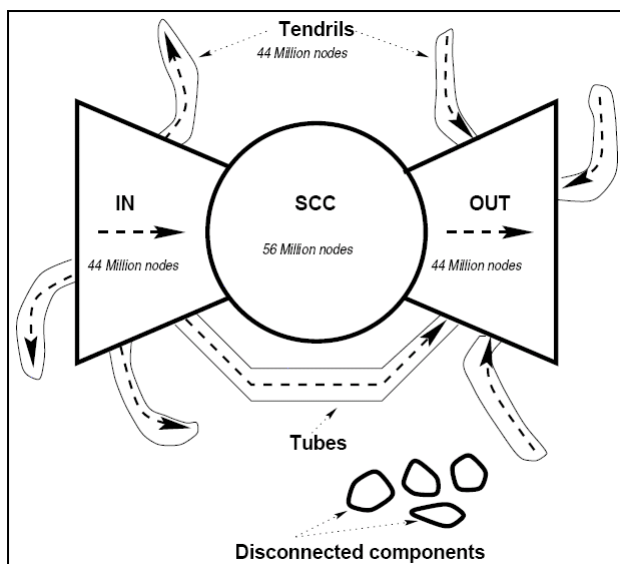
شکل ۲-۲ نتیجه مشابه را برای درجات خروجی نشان می‌دهد^۲. بهترین خط تطابق توزیع با $x = 2.72$ را برای ما می‌دهد. میانگین درجه خروجی برابر ۷.۲ می‌باشد. با استفاده از جستجوهای سطح اول و استفاده از گره‌های تصادفی Broder و سایرین توانستند که ساختار شکل ۲-۳ را برای گراف وب پیشنهاد کنند.^۳ این نتیجه از پیمایش Altavisto در ماه مه ۱۹۹۹ بدست آمده است که توسط پیمایش مشابه در همان سال جایگزین شده است. برای فهم این شکل ابتدا به بیان چند تعریف می‌پردازیم.

تعریف ۱. یک جزء همبند ضعیف مجموعه‌ای از صفحات است که به هم دسترسی دارند اگر پیوندها بدون توجه به جهت شان دنبال شوند.

وسیع ترین جزء همبند ضعیف در این پیمایش برابر ۱۸۶ میلیون گره بوده است که بیش از ۹۰ درصد از این پیمایش را شامل می‌شود.

^۲ اصل شکل از [۶]

^۳ اصل شکل از [۶]



نمایی از وب: 2-3

تعریف ۲. بطور مشابه یک جزء قویا همبند یک مجموعه از صفحات است که برای دو جفت صفحه (u, v) در این مجموعه یک مسیر جهت دار از u به v موجود است.

به زبان وب این بدان معناست که یک شخص از u می تواند توسط Hyper Link ها به v برسد. وسیع ترین جزء قویا همبند حدود ۵۶ میلیون گره دارد.

Broder و Kumar [۶] شکل پیشنهاد شده را Bowtie می نامد و ناحیه های مختلفی را در داخل Bowtie شناسایی می کند. هسته Bowtie که SCC نامیده است نمایشگر جزء قویا همبند گول آسای وب شامل ۵۶ میلیون گره است.

سمت چپ bowtie نمایشگر ۴۴ میلیون صفحه به نام IN است که نشاندهنده صفحاتی خارج از SCC است که به SCC راه دارند. پس یک کاوشگر وب با شروع از یک گره در IN می تواند به کلیه گره های در SCC دسترسی داشته باشد. IN بطور شهودی آن دسته از صفحات جدید است که خود به مقاصد جالب پیوند دارند ولی هنوز توسط هسته وب کشف نشده اند.

بطور مشابه مجموعه دیگر ۴۴ میلیونی از صفحات که سمت راست bowtie را نشان می دهند و OUT خوانده می شود که خصوصیت آن این است که هر صفحه ای از OUT توسط هر صفحه ای

در SCC قابل دسترسی است. بطور شهودی چنین صفحاتی صفحات مشهور و پیشکسوت اینترنت هستند که فقط به خودشان دسترسی دارند.

پس مدل ما توانایی کاوش bowtie را از چپ برآست نشان می دهد ولی عکس آن درست نیست.

نهایتا ناحیه چهارم که Tendrils خوانده می شود از صفحاتی تشکیل شده که به هسته راه ندارند و از هسته هم به آنها راهی نیست. این صفحات در حقیقت اجتماع عیوب مجموعه IN و OUT هستند. باز هم حدود ۴۴ میلیون از این گونه صفحات کشف شده است.

۴-۲ اندازه قطر وب

شکل Bowtie فاش کننده یک حقیقت غیرمنتظره در مورد همبندی وب است. برای اکثر صفحات v و u مسیری از u به v وجود ندارد. بطور دقیق تر اگر u در SCC اجتماع با IN و v در SCC اجتماع با OUT باشد در اینصورت مسیری بین u و v موجود است. احتمال قرار گرفتن u در IN و v در OUT حدود $\frac{1}{4}$ و احتمال قرار گرفتن v در OUT و u در IN حدود $\frac{1}{4}$ است پس احتمال این دو پیشامد مستقل حدود $\frac{1}{16}$ است.

پس در حدود ۷۵ درصد از صفحات وب به هم متصل نیستند.

آخرین تحلیل های Albert در [۲] قطر ۱۹ را برای وب پیش بینی کرده بود در حالی که اندازه گیری های واقعی تصویر متفاوتی به ما نشان می دهد اکثر جفت صفحات وب فاصله بی نهایت از هم دارند. به عنوان یک گزینه ما می توانیم میانگین فاصله اتصالی را تعریف کنیم بشرطی که بین دو صفحه مسیری واقع باشد. چنین میانگین ۱-۲ در جدول آمده است.

بدون جهت	خروجی	ورودی	نوع پیوند
6.83	16.18	16.12	میانگین فاصله اتصالی

میانگین فاصله اتصالی با توجه به نوع پیوند: 1-2

۵-۲ تحلیل صفحات میزبان

تا این جا ساختار گرافی وب را با واحدهای صفحات وب بررسی کردیم. در عمل وب یک ساختار سلسله مراتبی است با دامنه، صفحات میزبان و سایتها که نمایانگر در جات مختلف کنترل و مدیریت می باشند. برای درک بهتر ساختار وب ما باید کلان ساختار وب را در مقیاس پیوندهای بین وب سایت ها بررسی کنیم.

تعریف ۳. وب سایت: یک مجموعه از صفحات وب است که نمایانگر یک موجودیت در وب می باشد.

تعریف ۴. دامنه: متشکل از تعدادی میزبان وب است که جزء آخر مشترک در اسم خود دارند) بطور مثال com یا ir)

دامنه‌ها بیشتر وابسته به کشورها هستند با این وجود دامنه‌هایی هست که مرز جغرافیایی نمی شناسند. جدول ۲-۲ آماره در مورد سایز صفحات وب است.

	اکتبر 1999	2000 آگوست	2001 ژوئن
صفحات وب به میلیون	128.99	604.37	1,292
پیوندها به بیلیون	1.27	5.54	19.94

سایز وب: 2-2

	اکتبر 1999	2000 آگوست
بدون وزن	4.11	5.27
وزن دار	3.31	3.71

تخمین میانگین فاصله بین میزبان ها: 2-3

با توجه به تعریف ما از میانگین فاصله اتصالی در بخش قبل جدول ۲-۳ فاصله میانگین در جزء SCC نشان می دهد.

البته برای میزبان ها با توجه به تحلیل Henzinger در [۵] توزیع درجات در بین میزبان ها نیز از توزیع توانی یا Zipfian پیروی می کند که i مناسب برای این توزیع ۱.۶۲ برای درجات ورودی و ۱.۶۷ برای درجات خروجی است .

در نهایت مدلی که برای شبیه سازی گراف وب چه در مقیاس صفحه و چه در مقیاس میزبان ها پیشنهاد شده است به مدل نسخه برداری یا Copy Model مشهور است . بطور خلاصه برای ساختن چنین گرافی در هر گام با احتمال β ما یک گره از قبل موجود را انتخاب می کنیم و لینک های جدیدی را به آن طبق قانون زیر اضافه می کنیم :

- ابتدا یک گره تصادفی v از میان گره ها موجود و d تا لینک تصادفی آن را انتخاب می کنیم .
- در ادامه برای d تا از $i = 1, \dots, d$ امین لینک جدید u به یک گره تصادفی موجود با احتمال α اشاره می کند یا به مقصد i امین لینک v با احتمال $1 - \alpha$
- در نهایت با احتمال $1 - \beta$ یک گره جدید به گراف اضافه کرده و لینک های خروجی را بدان همانند مدل نسخه برداری اضافه می کنیم .
- این مدل نشان دهنده دو تغییر با گذشت زمان در وب است .

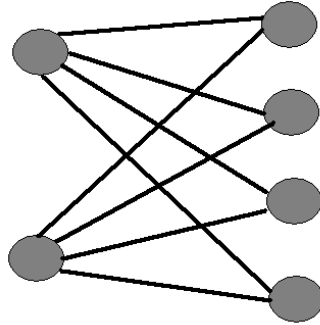
- (۱) امروز زمان گره های جدیدی را به گراف اضافه می کند یعنی صفحات یا میزبان های جدید.
- (۲) صفحات یا میزبان های کنونی می توانند پیوندهای جدیدی به گره های دیگر ایجاد کنند .

الگوریتم HITS و شاره در استخراج اجتماعات وب

۱-۳ الگوریتم HITS و هسته های دوبخشی

الگوریتم HITS بعنوان یک روش محلی در پیدا کردن اجتماعات وب نگاه می شود. یک روش کاملاً محلی آن روشی است که تنها خصوصیات یک همسایگی محلی حول و حوش دو گره را بررسی می کند تا تشخیص دهد که آیا این دو گره عضو یک اجتماع هستند. در مقابل روش های جهانی کلیه یال ها در گراف وب مورد بررسی قرار می دهد برای اینکه تصمیم بگیرد که آیا دو گره در داخل یک اجتماع قرار می گیرند یا خیر.

برای بررسی میزان شباهت یا ارتباط دو گره روش های سنتی و متریک های سنتی وجود دارد. اگر گراف وب را همراه با ماتریس مجاورت A در نظر بگیریم حاصلضرب $(AA^T)_{uv}$ و $(A^T A)_{uv}$ نشان دهنده میزان رجوع شونده و رجوع کننده برای جفت صفحه u و v می باشد. البته اینگونه متریک ها عملی رغم ظاهر تمیز و ساده‌ی خود مشکلات فراوانی در عمل خواهند داشت. اولاً هر دو ماتریس AA^T و $A^T A$ بسیار متراکم تر از خود A می باشد و این مشکلات ذخیره سازی و خواندن از حافظه را بدنبال می آورد. ثانیاً بعضی از صفحات ممکن است

نمای تصویری گراف دوبخشی $K_{2,3}$: 3-1

به (از) هزاران صفحه رجوع کنند (شوند). البته مشکل دوم با استفاده از نرمالسازی قابل حل می شود. پس تا این لحظه یک میزان برای تعیین شباهت میان دو صفحه با وجود تمامی مشکلات آن در پیاده سازی پیدا کرده ایم. ولی آنچه ما بدنبال آن هستیم حسی است که در دید کلی از اجتماعات وب بدست می آوریم.

تعریف ۵. یک گراف کامل دوبخشی یک گراف جهت دار است بصورتی که گره های آن قابل تقسیم به دو بخش L و R می باشند که $L \cup R = V$ و $L \cap R = \emptyset$ بدین صورت که از هر گره در L یالی به هر گره در R وجود دارد. ما چنین گرانی را به $K_{l,r}$ نشان می دهیم که $l = |L|$ و $r = |R|$. یک هسته دو بخشی یک زیرگراف از یک گراف است که اگر به تنهایی نگاه شود یک گراف کامل دوبخشی است. شکل ۱-۳ مثالی از $K_{2,3}$ می باشد.

الگوریتم HITS از Kleinberg در [۱۰] یک زیرگراف از گراف وب را می گیرد و دو وزن برای هر صفحه در این زیر گراف تولید می کند. این اوزان اغلب بعنوان امتیاز هاب و مرجعیت نامیده می شوند و ذاتا مربوط به خصوصیات ماتریسی آن بخش از گراف وب هستند که الگوریتم روی آن اجرا شده است.

بطور معنایی یک هاب یک صفحه از وب است که به تعداد زیادی از مراجع ارجاع می دهد و یک مرجع یک صفحه وب است که توسط تعداد زیادی از هاب ها ارجاع شده است. (به ماهیت بازگشتی این تعریف توجه کنید). دو امتیازی که HITS برای هر صفحه نشان می دهد نمایانگر این حقیقت است که تا چه حدی دو صفحه این دو خصوصیت را از خود نشان می دهد. اگر به شکل

۱-۳ مراجعه کنیم گره های سمت راست نمایانگر مراجع و سمت چپی ها نمایانگر هاب ها می باشند.

الگوریتم HITS در دو گام انجام می شود. مرحله اول پیش پردازشی است که معین می کند کدام زیرگراف از گراف غول آسای وب گزینش شود مادامیکه گام دوم یک پروسه عددی تکرار شونده است. راه معمول برای گزینش مجموعه دانه بصورت زیر است.

(۱) یک پرسش مورد علاقه را به یک موتور جستجو دلخواه بفرست.

(۲) مجموعه m صفحه اول جواب را انتخاب کن.

(۳) مجموعه را با کلیه همسایگان بلاواسطه یا با یک واسطه گسترش بده.

در انتخاب مجموعه دانه هیپورستیک های متفاوتی می تواند مورد استفاده قرار بگیرد. بعنوان مثال اینکه صفحات داخلی یک میزبان مورد استفاده قرار نگیرد یا اینکه یک حد بالایی برای تعداد صفحات گزینش شده بگذاریم و ...

با داشتن چنین مجموعه دانه ای اگر $G(V, E)$ این زیرگراف از گراف وب باشد و A ماتریس مجاورت G قسمت دوم تکرار شونده الگوریتم HITS دو بردار $1 \times |V|$ بعدی به نام های \mathbf{a}, \mathbf{h} را بصورت زیر بروز می کند. مقدار اولیه این دو بردار می تواند ۱ یا هر عدد ثابت دیگری باشد.

$$\mathbf{a} = A^T \mathbf{h} \quad (۱)$$

$$\mathbf{h} = A \mathbf{a} \quad (۲)$$

$$\mathbf{a} = \frac{\mathbf{a}}{\|\mathbf{a}\|^2}$$

$$\mathbf{h} = \frac{\mathbf{h}}{\|\mathbf{h}\|^2}$$

معادله ۱ بیانگر این مطلب است که امتیاز مرجع بودن یک صفحه جمع امتیازات هاب بودن صفحاتی است که او ارجاع داده اند. و مشابه معادله ۲ نشان دهنده این است که امتیاز هاب بودن یک صفحه جمع امتیاز مراجعی است که او به آنها ارجاع داده است.

۳-۱-۱ شباهت HITS و Pagerank

الگوریتم HITS از جهت تحلیل ماتریسی و نتایج بسیار شبیه الگوریتم Page Rank می باشد هر چند روش Page Rank در نگاه اول خیلی شبیه HITS نیست. به اختصار یک قدم زدن تصادفی در الگوریتم PageRank موجود است که در هر لحظه یکی از این دو عمل را انجام می دهد.

- با احتمال ε می تواند به یک صفحه در وب نقل مکان کند.
- با احتمال $1 - \varepsilon$ از روی یک یال تصادفی متصل به مکان فعلی اش به یک همسایه دیگر می رود.

در نهایت هر صفحه ای یک احتمال ثابت و قابل محاسبه ملاقات شدن توسط قدم زن تصادفی دارد که همان امتیازش است. بطور صریح محاسبه این امتیاز از طریق رابطه تکرار شونده زیر است.

$$r_v^t + 1 = \frac{\varepsilon}{n} + (1 - \varepsilon) \sum_{u:(u,v) \in E} \frac{r_u^t}{d_u^{out}} \quad (3)$$

kleinberg در [۱۰] با استفاده از الگوریتم خود توانست صفحاتی را که از نظرواژه‌ی کلیدی مشابه ولی معنای متفاوت بودند به رده های مناسب دسته بندی کند. بعنوان مثال صفحه حاوی لغت Jaguar به رده های مختلف مثل رده اتومبیل، حیوان، تیم فوتبال، بازی کامپیوتری، سیستم عامل و ... دسته بندی شد) با این مفهوم الگوریتم HITS قابل تبدیل به یک تشخیص دهنده اجتماع در وب است.

۳-۱-۲ سازگار شدن الگوریتم PageRank برای تشخیص اجتماعات

الگوریتم PageRank را نیز می توان براحتی برای تشخیص اجتماعات داخل وب سازگار کرد. معادله ی یکنواخت ۳ می تواند به یک معادله حساس به موضوع در روش PageRank تبدیل شود.

$$r_v^t + 1 = \varepsilon p_v + (1 - \varepsilon) \sum_{u:(u,v) \in E} \frac{r_u^t}{d_u^{out}} \quad (۴)$$

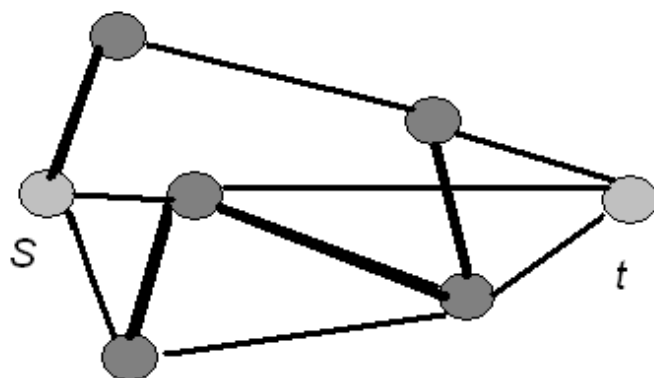
با این شرط که $\sum_v p_v = 1$ مقدار p_v برای تمامی v ها طوری انتخاب می شود که نشان دهنده احتمال پا گذاشتن قدم زن تصادفی به صفحه v با توجه به موضوع باشد. این پروسه می تواند در تشخیص اجتماعات بکار آید بدین ترتیب که برای دو صفحه یک امتیاز حساس به موضوع را محاسبه می کند. در نهایت صفحاتی که به اجتماع متعلق می باشند احتمال ملاقات شدنشان از یک حد بالاتر است و سایر صفحات خارج از اجتماع می باشند.

۲-۳ روش شماره بیشینه در بدست آوردن اجتماعات وب

روش شماره بیشینه از جمله روش هایی در بدست آوردن اجتماعات بحساب می رود که خصوصیات هر دو نوع روش های محلی و جهانی را دارد. این روش می تواند هم روی کل وب هم روی یک زیرگراف پیاده سازی شود. پیچیدگی این روش در بدترین حالت تابعی از سایز کل وب است با این وجود در عمل این پیچیدگی محدود به اندازه اجتماعی می شود که پیدا کرده است. همچنین اجتماعی که از این روش بدست می آید مطمئن تر و دقیق تر از سایر روش هاست. همانطور که انتظار می رود الگوریتم پیدا کردن اجتماع مساله را تبدیل به یک مساله پیدا کردن شماره بیشینه یا معادل آن برش کمینه می کند. در اینجا توضیح متخصری در مورد مساله شماره بیشینه و الگوریتم پیدا کردن آن می دهیم.

۱-۲-۳ مساله شماره بیشینه

یک شبکه آب رسانی را در نظر بگیرید. این شبکه از یک سری لوله های آب با قطرهای مختلف، یک سری گذرگاه و یک منبع و مقصد تشکیل شده است. قطرهای مختلف این لوله ها هر کدام ظرفیت شخصی برای گذردهی آب تعیین می کنند. شکل ۲-۳ نمونه ای از این شبکه را نشان می دهد.



یک شبکه شماره: 3-2

ما علاقمندیم که بدانیم این شبکه بیشترین مقدار آبی که از خود عبور می دهد چه اندازه است. بنابراین گراف $G(V, E)$ اگر معادل چنین شبکه ای باشد و برای کلیه یال های $(u, v) \in E$ نشان دهنده ظرفیت آن یال باشد و قرارداد می کنیم که $c(u, v) = 0$ در صورتی که بین u, v یالی موجود نباشد.

بنابراین با گرفتن دو راس s, t یک مساله شماره بیشینه $s - t$ پیدا کردن بیشترین مقدار شماره ایست که از راس s به راس t می تواند جریان داشته باشد بشرطی که هیچ یالی بیشتر از ظرفیت خود شماره تحمل نکند و هیچ شماره ای هم در بین راه نابود نشود. طبیعی است که آنچه این مقدار را تعیین می کند آن قسمت هایی از شبکه هستند که حکم گردی بطری را دارند. Fulkerman, Ford در تئوری «شماره بیشینه برش کمینه» خود ثابت کردند که بیشترین شماره ممکن در یک گراف با منبع s و مقصد t برابر با برش با کمترین ظرفیتی است که s را از t جدا می کند. الگوریتم چند جمله ای فراوانی برای حل مساله شماره بیشینه پیشنهاد شده است که بعنوان مثال می توان به [۱۱] اشاره کرد. در زیر یکی از ساده ترین این الگوریتم ها آمده است.

Max-Flow (graph $G(V, E)$; vertex: s, t)

set R = residual network of G

```

while R contains a directed path from s to t
  identify shortest augmenting path, P, from s to t
  set gamma = min (r(u, v) : (u, v) is in P)
  for all (u, v) in P
    set r(u,v) = r(u,v) - gamma
    set r(v,u) = r(v,u) + gamma
return R

```

این الگوریتم روی گراف باقیمانده ها کار می کند. این گراف اطلاعات یال ها را در مورد ظرفیت مصرفی و ظرفیت باقیمانده یا پتانسیل نگاه می دارد. در نهایت در هر دور الگوریتم یک مسیر روی گراف باقیمانده ها پیدا می کند که بتوان با ظرفیت های غیر صفر از s به t رسید. کمترین ظرفیت این مسیر به شماره کنونی اضافه می شود. اگر چنین مسیری وجود نداشت کار الگوریتم به پایان رسیده است.

۲-۲-۳ تعریف اجتماعات از دیدگاه شماره ای و الگوریتم تشخیص آن

تعریف ۶. یک اجتماع یک زیرگراف بارئوس $X \subseteq V$ می باشد به نحوی که برای هر گره $v \in X$ تعداد یال هایی که v را به گره های داخل X متصل می کند کمتر از تعداد یال هایی که v را به خارج وصل می کند نباشد.

توجه کنید که چنین تعریفی قدری بازگشتی است چون دارای غالبی مشابه غالب زیر است که یک ایرانی بیشتر ایرانیان را می شناسد تا غیر ایرانیان را. البته تشخیص فراگیر و دقیق اجتماعات با این تعریف یک مساله NP تمام است با این حال روش شماره بیشینه اکثر اجتماعات قوی را باز می گرداند. ورودی الگوریتم تشخیص اجتماع گراف G و یک مجموعه دانه S و یک پارامتر α می باشد. رویه، یک گراف از روی G به انضمام دو گره منبع و مقصد مصنوعی می سازد از منبع با ظرفیت نامحدود به کلیه رئوس مجموعه دانه و از مقصد با ظرفیت تعیین شده توسط α به کلیه رئوس کل

گراف متصل می‌کنیم. پارامتر α در حقیقت یک پارامتر کوچک کردن می‌باشد که کار آن تنظیم کردن سائز اجتماعات استخراجی می‌شود. بعد از ساخت رویه ما الگوریتم شماره پیشینه را به عنوان یک زیر برنامه فراخوانی می‌کند و در انتها آن بخش از گراف باقیمانده‌ها که در مولفه همبندی S قرار می‌گیرد به عنوان یک اجتماع منطبق با تعریف؟ برگردانیده خواهد شد.

در زیر شبکه‌کد الگوریتم فوق را آمده است. دو الگوریتم مرتبط به نام‌های Approximation-com و com-cluster داریم که هر دو آنها الگوریتم com را به عنوان یک زیر برنامه فراخوانی می‌کنند. الگوریتم com هنگامی است که کل گراف در حافظه اصلی گنجانیده شود.

```
com ( graph G(V,E), S, alpha)
  set V_alpha = V union {s, t}
  set E_alpha = E union {(v, t): v in V} union {(s, u) : u in S}
  set c(v, t) = alpha, for all v in V
  set c(s, u) = not_bounded for all u in S
  let R be the output of Max-Flow(G_alpha, s, t)
  set X be the member of smallest connected component about s in R
  return X - {s}
```

همانطور که می‌دانیم گراف غول‌آسای وب دارای چنین خصوصیتی نمی‌باشد. پروسه Approximation-com هنگامی به کار می‌رود که تنها بخش کوچکی از گراف بتواند در حافظه نگاه دار می‌شود. این الگوریتم از یک خزش^۱ اندازه ثابت برای یافتن یک اجتماع تقریبی استفاده می‌کند و قوی‌ترین اعضای این اجتماع را به عنوان دانه‌های مرحله بعدی در نظر می‌گیرد.

Approximation-com (S, k)

```
while number of iteration is less than an optional number
  set G to a crawl from S of depth k
```

^۱crawl

```

let alpha be the size of S
set X = com(G, S, alpha)
rank all v in V by number of edges in X
add highest ranked non-seed vertices in X to S
return X

```

رویه com-cluster برای پیدا کردن کلیه اجتماعات در یک گراف بکار می رود. این الگوریتم هم فرض می کند که کلیه گراف در داخل حافظه اصلی جا می گیرد.

```

com-cluster(graph G(V,E), alpha)
set S = V
while exists s in S
  set X = com(G, {s}, alpha)
  for all v in V
    set cluster(v) = s
  set S = S - X
return cluster labels of v in V

```

۳-۲-۳ تحلیل

بخاطر آورید که ما از پارامتر به عنوان پارامتر کوچک کردن نام بردیم. قضیه زیر در مورد پارامتر صادق است.

قضیه ۷. اگر X اجتماعی باشد که توسط الگوریتم com با پارامتر α پیدا شده است برای هر P و Q به نحوی که P و Q را افزاز کنند حدود زیر برقرار است.

$$\frac{f(x, V - X)}{|V - X|} \leq \alpha \leq \frac{f(P, Q)}{\min(|P|, |Q|)}$$

اثبات قضیه فوق در [۱۴] یافت می شود. در حقیقت α یک حد بالایی برای ظرفیت های برون اجتماعی و یک حد پایینی برای ظرفیت های درون اجتماعی است. بنابراین الگوریتم com تضمین می کند که اعضای اجتماع بطور متراکم به هم مرتبط و یال های کمی به خارج داشته باشند. در دو حالت اکسترم انتخاب یک مقدار کوچک برای α مثلاً نزدیک به صفر منتهی به پیدا کردن کل گراف به عنوان یک اجتماع می شود در صورتی که انتخاب یک مقدار بزرگ برای α منجر برگرداندن n اجتماع تک عضوی^۲ خواهد شد که در اینجا منظور از یک مقدار بزرگ می تواند عددی بزرگتر از مجموع ظرفیت های گراف باشد. واقعیت دیگر در مورد الگوریتم Com-Cluster است. برای دو اجتماع تشخیص داده شده X_u و X_v در این الگوریتم یکی از سه حالت زیر برقرار است.

$$X_v \subseteq X_u \text{ یا } X_u \subseteq X_v \text{ یا } X_u \cap X_v = \emptyset$$

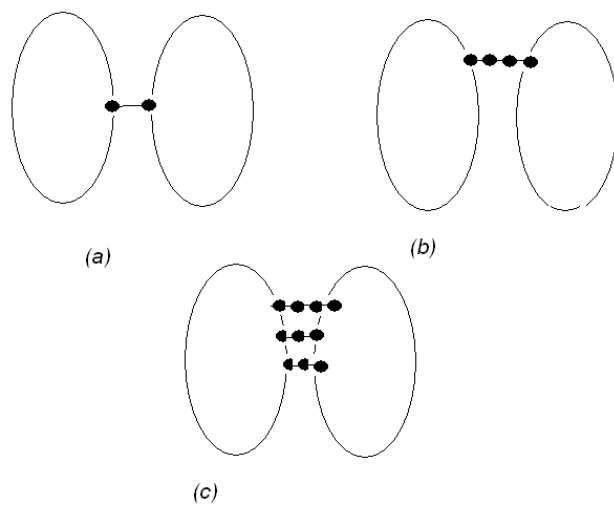
۳-۳ بحث و مقایسه

در الگوریتم درخت های برش کمینه در فصل آتی مفهوم کانون مطرح می شود. مقصد مصنوعی که ما در این الگوریتم برای گراف انتخاب کردیم در رابطه تنگاتنگی با کانون درخت برش کمینه قرار می گیرد. به فصل مربوطه مراجعه کنید.

در مقایسه با الگوریتم HITS باید گفت یک مزیت الگوریتم های مبتنی بر شماره نسبت به الگوریتم HITS و مشابهها Pagerank عدم حساسیت الگوریتم شماره به فواصل می باشد. در شکل ۳-۳ این تفاوت به تصویر کشیده شده است. در حالی که اشکال a و b از دید الگوریتم شماره یکسان است. ولی یک قدم زن تصادفی بسیار بین این دو شکل تمییز قائل می شود چون در شکل b مسیر طولانی باید برای رسیدن به یک جزء از جزء دیگر طی شود در حالی که شکل a چنین خصوصیتی ندارد.

در مقابل شکل c با وجودی که از نظر یک قدم زن تصادفی با شکل a توفیری ندارد ولی الگوریتم شماره این دو شکل را بسیار متفاوت می یابد زیرا شماره های پیشینه در این شکل با هم تفاوت می کنند.

^۲ singleton



چند گراف مشابه با خصوصیات برشی و ماتریسی گوناگون: 3-3

فصل ۴

الگوریتم‌های دسته‌بندی با درخت‌های برش کمینه

۴-۱ درخت‌های برش کمینه

در این بخش نگاه اجمالی به تئوری درخت‌های کمینه می‌کنیم. حتی الامکان از پرداختن به اثبات‌ها در بیان تئوری‌ها خودداری می‌کنیم. این بخش زمینه را برای معرفی روش Tsioutsioulis در دسته‌بندی فراهم می‌کند.

بیان درخت برش کمینه در وهله اول توسط Hu, Gomory در [۱۲] مطرح شد. از این بابت این داده ساختار بر درخت Gomory-Hu نیز خوانده می‌شود. فرض کنید $G(V, E)$ گراف بدون جهت و $V = \{v_1, \dots, v_n\}$ مجموعه گره‌ها و $f(i, j)$ یا f_{ij} شماره بیشینه بین دو گره $v_i, v_j \in V$ باشد.

قضیه ۸. مجموعه اعداد نامنفی $f_{ij} = f_{ji}$ که در آن $i, j \in \{1, \dots, n\}$ یک بیشینه شماره معتبر است اگر و فقط اگر

$$f_{ik} \geq \min\{f_{ij}, f_{jk}\}$$

برای تمام $i, j, k \in \{1, \dots, n\}$

نتیجه: با استفاده از استقرا می‌توان گفت برای هر دنباله v_1, \dots, v_n از گره‌ها

$$f_{1n} \geq \min\{f_{12}, f_{23}, \dots, f_{(n-1)n}\}$$

علاوه بر این برای هر شبکه بدون جهت G یک درخت T وجود دارد که حاوی شماره‌های بیشینه آن شبکه است.

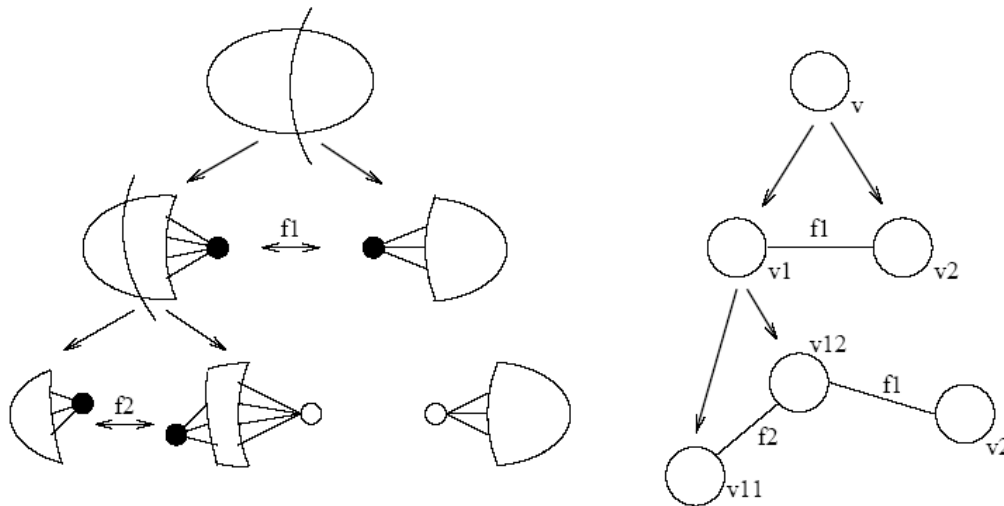
قضیه ۹. اندازه‌ی شماره بیشینه بین دو گره v_i و v_j از گراف G برابر با کمترین مقدار مجموعه $\{f_{ia}, f_{ab}, \dots, f_{dj}\}$ است که $v_{ia}, v_{ab}, \dots, v_{dj}$ وزن یال‌هایی هستند که مسیر یگانه بین دو گره v_i و v_j را در T تشکیل می‌دهند.

به T درخت معادل شماره بیشینه نیز می‌گویند. از بین درخت‌های با چنین خصوصیت حداقل یک درخت وجود دارد که حذف یال با کمترین وزن روی مسیر بین v_i و v_j به ما برش کمینه بین v_i و v_j را در G می‌دهد. به چنین درختی با این خصوصیت درخت Gomory-Hu، درخت برش کمینه یا درخت برش می‌گویند. با در دست داشتن چنین درختی (که باز هم تاکید می‌کنم می‌تواند یگانه نباشد) ما شماره کمینه بین دو گره v_i و v_j را در گراف G با پیمودن مسیر بین v_i و v_j (که می‌دانیم چنین مسیر یگانه‌ای وجود دارد) و یافتن یال کمینه e_{\min} پیدا می‌کنیم.

۴-۱-۱ الگوریتم ساختن درخت

نسخه‌ی اصلی

- قدم اول: الگوریتم شماره بیشینه را بین دو گره در گراف اجرا کن
 - قدم دوم: خروجی قدم اول بعنوان وزن یال بکار رفته در درخت برش بکار می‌رود.
- این الگوریتم بعد از ساختن $n-1$ یال به اتمام می‌رسد. در مرحله یک چشمه و چاهک جدید برای قدم اول انتخاب می‌شود.



الگوریتم ساختن درخت برش کمینه: 4-1

با این تفاوت که زیرمجموعه‌های شخصی در گراف اصلی بصورت یک گره نگاه می‌شوند. این کار باعث می‌شود که ورودی مساله ما کوچک‌تر از گراف اصلی شود و در نهایت مرتبه الگوریتم کاهش یابد. این دو مرحله بطور روشن در شکل ۴-۱ به تصویر کشیده شده‌اند.^۱

با گرفتن گراف اولیه به عنوان ورودی یک گره از درخت برش ایجاد می‌کنیم. آنگاه دو گره دلخواه از گراف انتخاب کرده و الگوریتم برش کمینه را روی آن دو اجرا می‌کنیم. با محاسبه برش کمینه درخت اولیه برش ما که یک گره داشت به دو گره جدید با یک یال که وزن آن برابر مقدار برش کمینه حاصله است تبدیل می‌شود. این دو گره نماینده دو بخش در دو سوی برش کمینه حاصله در گراف اصلی می‌باشد در دور بعدی الگوریتم ما روی یکی از این دو بخش برش کار می‌کنیم. این گراف جدید همان گراف اولیه است با این تفاوت که بخش دیگر برش را در یک گره منقبض کردیم و طبیعتاً حالا باید دو گره‌ای که برای الگوریتم برش کمینه انتخاب می‌شوند شامل این گره انقباضی نباشند. حال شماره پیشینه جدید به ما دومین یال را از این درخت حاصل می‌کند. باز بطور مشابه گره‌ای که نماینده یک بخش از برش در گراف اولیه ما بود به دو گره جدید تقسیم و

^۱(اصل شکل از [۱۳])

یال بین این دو گره وزن شماره بیشینه را خواهد گرفت. اکنون گره سابق این درخت را باید به یکی از این دو گره جدید اتصال دهیم. طبیعتاً اتصال ما به آن بخش از برش خواهد بود که در گراف واقعی هم گره ما بدان اتصال دارد.

با اجرای الگوریتم بدین منوال نهایتاً به مرحله‌ای می‌رسیم که هر گره ما در درخت برش متناظر با زیرگرافی از گراف اصلی است که دقیقاً یک گره اصلی دارد. این اتفاق دقیقاً بعد از $n-1$ دور از الگوریتم می‌افتد. البته Gomory و Hu نشان دادند که این درخت براساسی همان درخت برش کینه با همان خصوصیات دلخواه است.

نسخه Gusfield

نسخه‌های بهتری نیز از این الگوریتم ارائه شده است. در الگوریتم Gusfield در [۱۳] دیگر احتیاجی به انقباض یک گره وجود ندارد. با وجودی که این الگوریتم هم در $n-1$ محاسبه شماره بیشینه انجام می‌شود ولی تمام اعمال روی همان گراف اصلی انجام خواهد شد. در این حالت درخت برش ما با یک گراف ستاره آغاز می‌شود و در مراحل اجرای الگوریتم یال‌ها با توجه به مقادیر محاسبه شماره بیشینه بین گره‌ها جابجا خواهند شد. در زیر این الگوریتم آورده شده است.

`GusfieldCutTree(G(V, E))`

Initialize:

Number all nodes in V from 1 to n

Let p be an n length vector initialized to 1

// p corresponds to star tree T ,

// such that every node s points to $p[s]$

For $s=2$ to n do

Compute a minimum cut between nodes s and $t = p[s]$ in G

Let X be the set of nodes on the s side of the cut

```

Let f(s,t) be the value of the minimum cut found
Set fl[s] = f(s, t)
For i = 1 to n do
  If ((i != s) and ( i in X) and (p[i] == t))
    Set p[i] = s
If (p[t] in X)
  Set p[s] = p[t]
  Set p[t] = s
  Set fl[s] = fl[t]
  Set fl[t] = f(s,t)

Output p and fl /* p corresponds to the minimum cut tree T
with edges defined between all nodes i and p[i];
the weights of edge (i, p[i]) is equal to fl[i].

```

به رغم ساختار ساده تر الگوریتم Gusfield و پیچیدگی محاسباتی برابر این الگوریتم با الگوریتم اولیه Gomory-Hu نتایج تجربی در [۱۳] نشان می‌دهد که الگوریتم دوم با استفاده از هیوریستیک‌های مناسب قوی تر کار می‌کند.

به علت اعمال انقباضی پیاده سازی بهینه برای الگوریتم Gomory-Hu به بداهت الگوریتم Gusfield نمی‌باشد. یک پیاده سازی ناظریف حافظه جدیدی به گره انقباضی و یالش اختصاص می‌دهد. ممکن است که $\Omega(n^2)$ اختصاص حافظه حتی برای یک گراف خلوت بیانجامد. Tsioutsoulis در [۱۳] پیاده سازی پیشنهاد کرد که از مرتبه $O(\lg n)$ رکورد برای گره‌های اضافی و $O(n)$ رکورد برای یال‌های اضافی بود.

۲-۴ الگوریتم‌های دسته‌بندی و استخراج اجتماعات

یک درخت برش حاوی اطلاعات ارزشمندی هم درباره برش‌های کمینه در یک گراف G و هم ساختار درونی G می‌باشد.

لم ۱۰. فرض کنید u, v دو گره از درخت برش T باشد و $P(u, v)$ مسیر (واحد) بین آن دو در T . آنگاه برای هر گره w روی این مسیر شماره بیشینه $f(u, w)$ کمتر از $f(u, v)$ نخواهد بود.

لم ۱۱. فرض کنید $u, v, P(u, v)$ به ترتیب دو گره و مسیر بین آن دو در درخت برش کمینه باشد و $p = P(u, v)$ اندازه یا طول این مسیر برابر با تعداد برش‌های کمینه یگانه s, t (متفاوت) در گراف است.

بطور کلی در درخت برش کمینه گره‌های با اتصال شدید به هم کنار یکدیگر قرار می‌گیرند. البته چنین حقیقتی اگر بخواهد برای کلیه گره‌ها صدق کند در تعارض با همسایه‌های ایجاد شده قرار می‌گیرد و ممکن است منتهی به ساختار ضعیف درختی شود. در زیر الگوریتم پایه دسته‌بندی را می‌آوریم و توضیح می‌دهیم که براساس ساختار درخت برش دسته‌بندی می‌کند.

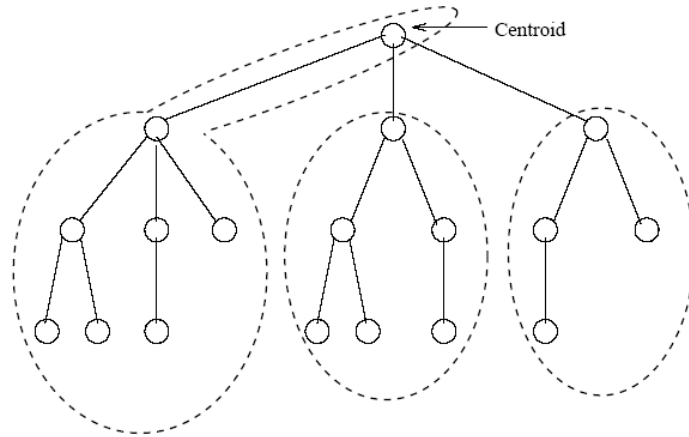
۱-۲-۴ الگوریتم

تعریف کانون (centroid) درخت برش :

تعریف ۱۲. کانون درخت T با مجموعه گره‌های V به صورت زیر تعریف می‌شود:
اگر $v \in V$ و T_1, T_2, \dots, T_d اجزا همبند ایجاد شده با حذف v از T باشند در اینصورت اگر $|T_i|$ اندازه زیردرخت T_i تعریف شود آنگاه تعریف می‌کنیم برای هر گره:

$$N(v) = \max\{|T_i|\}$$

که در آن $1 \leq i \leq d$



تشخیص کانون در درخت: 4-2

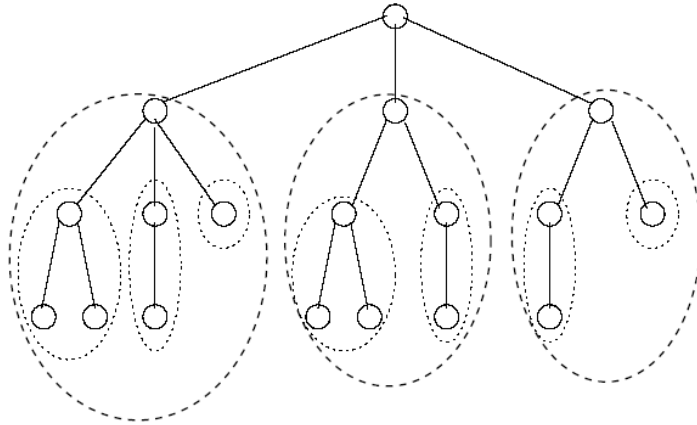
حال کانون درخت T تعریف می‌شود ([50]) به عنوان آن گرهی از درخت که می‌نیمم می‌کند مقدار $N(v)$ را روی تمام گره‌ها در V یا

$$N(v_c) = \min\{N(v)\}$$

که در آن $v \in V$

بطور شهودی آنچه از کانون یک درخت انتظار می‌رود آن است که با حذف کانون از درخت اجزا همبند باقیمانده بصورت حتی الامکان عادلانه تقسیم شده باشند. درخت T می‌تواند تا دو کانون داشته باشد که در اینصورت انتخاب هر کدام در الگوریتمی که در ادامه توضیح می‌دهیم امکان پذیر است. پس اگر $G(V, E)$ گراف ما برای عمل دسته‌بندی باشد و T درخت برش کمینه مربوط به آن باشد الگوریتم در ابتدا با پیدا کردن کانون درخت شروع می‌کند. با تشخیص کانون ما درخت را از کانون آویزان می‌کنیم حال هر زیردرخت زیر کانون که هم اکنون ریشه درخت است یک دسته (cluster) است که توسط الگوریتم برگردانده می‌شود. حال باید ریشه را به یکی از این دسته‌ها منصوب کنیم. اگر e_c یال با بیشترین وزن متصل به v_c باشد ما v_c را به آن زیردرخت منسوب می‌کنیم که یک سر e_c در آن است.

شکل ۴-۲ مثالی از الگوریتم ماست.



دسته‌بندی سلسله‌مراتبی (اصل شکل از [۱۳]): 4-3

۲-۲-۴ دسته‌بندی سلسله‌مراتبی

این الگوریتم به ما قدرت دسته‌بندی سلسله‌مراتبی را نیز فراهم می‌کند بدین ترتیب که می‌توانیم هر زیردرخت را خود به عنوان یک مجموعه جدید برای دسته‌بندی بگیریم. فقط بخاطر داشته باشید که دیگر لازم نیست کانون برای این زیردرختان تعیین کنید. شکل ۳-۴ یک نمونه از دسته‌بندی سلسله‌مراتبی را نشان می‌دهد.^۲

۳-۲-۴ نرمالسازی یالی

همانطور که Benezur [۱۵] نشان داد تعریف درخت برش کمینه فقط درباره گراف‌های بی‌جهت است. برای تصمیم‌گیری این تعریف به گراف‌های جهت‌دار که موضوع مورد علاقه ما از جهت صفحات وب است از نرمالسازی یالی خروجی استفاده می‌کنیم.

اگر $G(V, E)$ گراف جهت‌دار و $v \in V$ یکی از گره‌های این گراف باشد آنگاه فرض کنید $E_{in}, E_{out} \subseteq E$ به ترتیب مجموعه یالی ورودی و خروجی باشد و برای هر یالی $w(x, y)$ که $(x, y) \in E$ وزن مربوطه باشد.

^۲(اصل شکل از [۱۳])

حال نرمالسازی یال خروجی $(u, v) \in E_{out}$ بصورت زیر انجام می‌شود.
 یال (u, v) را بدون جهت فرض کرد و وزن جدید $\hat{w}(u, v)$ را به صورت زیر به آن نسبت می‌دهیم

$$\hat{w}(u, v) = \frac{w(u, v)}{\sum_{x \in V} w(u, x)}$$

لم ۱۳. این نحوه نرمالسازی G را به یک گراف بدون جهت با یال‌های نرمال شده تبدیل می‌کند.
 اثبات. اثبات. هر گرهی وزن‌های یال‌های خروجی خود را فقط نرمال می‌کند و چون هر یال خروجی یک و فقط یک گره است حکم ثابت می‌شود. چنین نحوه نرمالسازی ما را به یاد روش‌های قدیمی‌تر مانند HITS, Pagerank می‌اندازد. □

۴-۲-۴ کاربرد در WWW

حال آماده ایم که الگوریتم پیشنهادی را روی وب اجرا کنیم. برای این کار یک مجموعه اولیه دانه از صفحات وب که مربوط به کلمه کاوشی (Search term) ما هستند پیدا می‌کنیم. این کار می‌تواند از طریق یکی از موتورهای جستجو انجام گیرد. حال برای ساختن گراف ورودی مساله ما به تمام بیوندهای ورودی و خروجی مجموعه دانه نگاه می‌کنیم. گره‌های این گراف مجموعه دانه بعلاوه همسایگان بلاواسطه‌ی آنها می‌باشد. حال می‌توان با نرمالیزه کردن و اجرای الگوریتم قسمت ۴-۲-۱ دسته‌ها را محاسبه کرد.

در اینجا خلاصه‌ای از نتایج تجربی که Tsioutsoulis در [۱۳] در مقایسه با Kleinberg برای عبارات Java, censorship, search engine, gates بدست آورد می‌آوریم. هر دو از یک مجموعه دانه با همان استراتژی همسایه با یک فاصله (بلاواسطه) استفاده کردند. بعنوان مثال دسته‌هایی که Kleinberg برای عبارت search engine پیدا کرده بود ۵ تای اولش همان موتورهای جستجوی معروف بودند در حالیکه Tsioutsoulis در ۵ تای اول و عموماً ۱۰ الی ۲۰ تای اول صفحاتی در مورد موتورهای جستجو و فراموتورهای جستجویی که قادر به پشتیبانی جستجوهای چندگانه و همزمان هستند بود. بنابراین نتیجه الگوریتم حداقل در این مثال بهتر و کلی‌تر از آنچه Kleinberg با الگوریتم خود پیدا کرده بود بدست آمد.

۵-۲-۴ الگوریتم دسته‌بندی اجتماعات

بهتر است در ابتدا یک مساله بسیار مشابه که NP-complete می‌باشد را مطرح کنیم. این الگوریتم به دسته‌بندی k -برشی کمینه یا (Minimum k -cut clustering) معروف است و بصورت زیر بیان می‌شود.

اگر $G(V, E)$ یک گراف با مجموعه گره V و یال E باشد آنگاه طوری این گراف را به k زیرمجموعه افزایش دهید که جمع یال‌های بین زیرمجموعه‌ها کمینه در بین تمام افزایش‌ها باشد. با استفاده از درخت برش کمینه می‌توان یک راه حل تقریبی با ضریب تقریب $(2 - \frac{2}{k})$ برای این مساله ارائه کرد. بدین صورت که اگر T درخت برش گراف بدون جهت $G(V, E)$ باشد با برداشتن $1 - k$ یال با کمترین وزن در T جزء همبند بدست می‌آید که G را به k دسته افزایش می‌کند. مجموعه وزن یال‌های بین این دسته‌ها یک تقریب $(2 - \frac{2}{k})$ از مساله دسته‌بندی k -برش کمینه است. با استفاده از این روش Wu, Lealy در [۱۶] ثابت کردند که هر برش درون دسته‌ای حداقل به بزرگی هر برش بین دسته‌ای در G می‌باشد. با وجود اینکه این روش بسیار راحت و سراسر است می‌باشد ولی عیب اساسی دارد. با توجه به این حقیقت که معمولاً یال‌های سبک در درخت برش ما به سمت برگ‌ها می‌روند (به تعریف کانون در بخش ۴-۲-۱ مراجعه کنید) دسته‌بندی با استفاده از این روش دسته‌های بسیار نامتعادلی تولید می‌کند که بیشتر آنها گره‌های تکی هستند. این همان نتیجه ایست که tsioutsoulis روی مجموعه داده هایش نیز بدست آورد.

۶-۲-۴ اجتماع ازوتریک

حال آماده می‌شویم که منظور خود را از یک اجتماع در وب بیان کنیم.

تعریف ۱۴. یک اجتماع (اجتماع ازوتریک) یک مجموعه S است بنحوی که بازای هر $s \in S$

$$\sum_{v \in S} w(s, v) \geq \sum_{v \in V-S} w(s, v)$$

به بیان دیگر یک عضو اجتماع بیشتر به اجتماع خود تعلق دارد تا به بیگانگان!!
در نهایت ما به دنبال دسته‌هایی هستیم که شبکه ما را به اجتماعات ازوتریک افزایش کند.
Tsioutsoulis در [۱۳] نشان می‌دهد که این مساله در حالت عمومی NP-تمام می‌باشد (کاهش
به مساله افزایش متعادل) با این وجود برای بسیاری از درختان برش می‌توان جواب‌های خوبی از این
مساله یافت.

۷-۲-۴ تعمیم تعریف اجتماع Esoteric

S یک اجتماع p -esoteric است هرگاه به ازای هر $s \in S$

$$\sum_{v \in S} w(s, v) \geq p \sum_{v \in V} w(s, v)$$

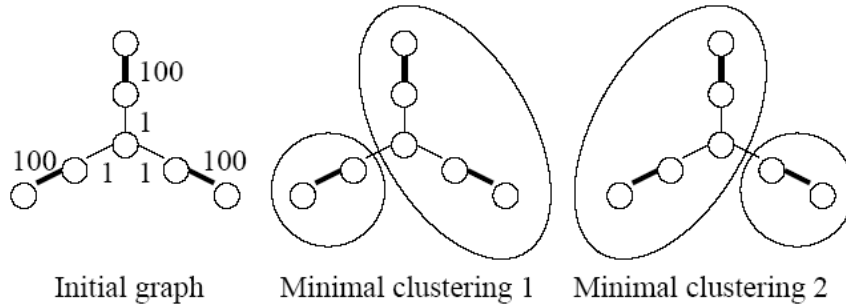
یا معادلا

$$(1-p) \sum_{v \in S} w(s, v) \geq p \sum_{v \in V-S} w(s, v)$$

دقت کنید که با چنین تعریف تعریف قبلی ما از اجتماع $1/2$ -esoteric است. به پارامتر p که
همواره از $1/2$ کوچکتر پارامتر همبندی گفته می‌شود.

خصوصیات

- (۱) هر اجتماع باید حداقل دو گره داشته باشد. این ادعا نتیجه مستقیم تعریف اجتماع است.
- (۲) اگر C یک دسته بندی G باشد و \hat{C} از اجتماع اجتماعات C تشکیل شده باشد آنگاه \hat{C} هم
یک دسته بندی مجاز برای G است. علت این خصوصیت نیز واضح است زیرا عمل اجتماع
فقط وزن های داخل یک اجتماع را افزایش می‌دهد.
- (۳) اگر یک رده بندی غیرقابل کاهش را آن طور تعریف کنیم که رده بندی است که نتوان در آن
دسته‌ها را کوچکتر کرد آنگاه می‌توان بجای یک دسته بندی کمینه چندین دسته بندی



دو دسته‌بندی غیر قابل کاهش متفاوت (شکل از [۱۳]: 4-4)

غیر قابل کاهش داریم. شکل ۴-۴ چنین دسته‌بندی را نشان می‌دهد.^۳

۳-۴ جستجو برای اجتماعات

برای یافتن اجتماعات ازوتریک از طریق شاره بیشینه باید پارامتر همبندی را برابر $1/2$ قرار دهیم. اگر t, s به ترتیب منبع و چاهک برای گراف G باشد برش کمینه $s - t$ را در G در نظر بگیریم. اگر S, T به ترتیب زیر مجموعه‌های دو سوی این برش باشد بنحوی که $s \in S$ و $t \in T$ آنگاه بردن هر گره از $S - s$ به T اندازه برش را فقط می‌تواند کاهش دهد. این حقیقت برای عکس قضیه نیز درست است. پس هر گره در S, T بیشتر به خودشان متصل هستند تا قسمت متقابل پس S, T هر دو از قانون اجتماع پیروی می‌کنند. استثنا می‌تواند خود t, s باشد. براساس چنین مشاهده‌ای ما می‌توانیم با استفاده از درخت برش کمینه اجتماعات ازوتریک را بیابیم.

تعریف ۱۵. یک یال e در درخت برش کمینه T یک یال کمینه محلی است اگر هر دو سر v_1 و v_2 از e حداقل یک یال در T داشته باشند که وزنش بزرگتر از e باشد.

^۳ اصل شکل از [۱۳]

برای پیدا کردن چنین یال‌هایی یک راه ساده وجود دارد. برای تمامی گره‌های داخل درخت یال با بیشترین وزن مجاورشان را علامت بزنییم. حال تمام یال‌های بدون علامت یال کمینه محلی هستند.

لم ۱۶. اگر T یک درخت برش کمینه برای گراف G باشد و e یک یال کمینه محلی آنگاه حذف e دو اجتماع از و تریک به ما می‌دهد.

اثبات. برای اثبات فرض کنید v_1 و v_2 دو گره دو سر e می‌باشند. یال e معادل برش کمینه v_1, v_2 در G می‌باشد. بنابراین دو طرف این برش دو اجتماع از و تریک هستند اگر

(الف) هیچ کدام از این دو طرف یک گره تنها نباشند.

(ب) اگر انتقال هر کدام از v_1 و v_2 طرف دیگر یک برش با اندازه کمتر به ما ندهد.

الف هیچ‌گاه اتفاق نمی‌افتد چون در اینصورت e به عنوان پروزن ترین یال v_1 یا v_2 از کاندیداهای یال‌های کمینه محلی خط می‌خورد.

(ب) فرض کنید بدون از دست دادن کلیت که انتقال v_1 به قسمت برش v_2 با مقدار کمتری به ما بدهد. اگر $e_{\max} \neq e$ یال مجاور v_1 با بزرگترین وزن باشد و v_{\max} و v_1 دو سر این یال باشند با

انتقال v_1 به قسمت دیگر برش v_1 از v_{\max} جدا می‌شود پس نمی‌تواند وزنی کمتر از e_{\max} پیدا کند که می‌دانیم بیشتر از وزن e می‌باشد پس ثابت شد که حذف e اجتماعات از و تریک معتبر برای

ما پدید می‌آورد بر اساس لم قبل با الگوریتمی ارائه می‌کنیم که اجتماعات از و تریک را در G پیدا می‌کند. □ ایده اصلی الگوریتم بدین صورت است که ابتدا دو اجتماع پیدا می‌کنیم و یکی از این دو را

به یک گره منقبض می‌کنیم بعنوان مثال X در تکرارهای بعدی باز دو اجتماع دیگر پدید می‌آید و هر دفعه یکی از این اجتماعات جذب X می‌شود. در نهایت الگوریتم هنگامی خاتمه می‌یابد که

هیچ یال کمینه محلی در درخت برش کمینه ما وجود نداشته باشد.

در زیر الگوریتم شبکه مربوطه آمده است.

Community Cluster ($G(V,E)$)

Construct T as the minimum cut tree of G

if (T has no local minimum edge)

Return single, trivial cluster covering entire G

```

else
  let X={ } be empty set
  let e be a local minimum edge in T
  while (e non-empty)
    let S1, S2 be the two communities formed after removing e
    if (X= { })
      output S1, S2 as next community and
      contract if into X
    else if (X is a member of S1)
      output Sz as next community and contract it
      into X
    else { // if X S2
      Output S1 as next community and contract it
      into X}
  Calculate new minimum cut tree T
  Let e be a local minimum edge in new T
Return

```

Tsioutsoulis در [۱۳] نشان داد که این الگوریتم اجتماعات ازوتریک $1/2$ تولید می‌کند. توجه کنید که این الگوریتم لروماگراف را افزاز نمی‌کند بلکه تنها اجتماعات را از G استخراج می‌کند.

۱-۳-۴ تسهیل شرایط

همانطور که توضیح داده شد الگوریتم در مراحل مختلف خود یال‌های کمینه محلی جدید را از روی درخت محاسبه شده پیدا می‌کند. یک ایده جایگزین این است که در همان ابتدا تمام یال‌های کمینه محلی را از درخت اولیه برش حذف کنیم. البته دسته‌هایی که با این عمل حاصل می‌شوند

ممکن است لزوماً اجتماعات از و تریک نباشند ولی با این حال می‌توانیم لم زیر را برای آنها ثابت کنیم.

لم ۱۷. اگر T درخت کمینه برش از گراف G باشد و C دسته بندی حاصله توسط حذف کلیه یال‌های کمینه محلی از T باشد آنگاه فرض کنید v_1 و v_2 دو گره در یک دسته مشابه از C باشد و $c(v_1)$ و $c(v_2)$ بزرگترین مقدار برش کمینه که v_1 و v_2 را از سایر گره‌ها در G جدا می‌کند باشد آنگاه برای برش کمینه $c(v_1, v_2)$ بین v_1 و v_2 ما داریم

$$c(v_1, v_2) \geq \min(c(v_1), c(v_2))$$

اثبات. می‌دانیم گره‌های روی مسیر بین v_1 و v_2 متعلق به C_1 یعنی دسته ایست که v_1 و v_2 بدان تعلق دارند. نظر به وزن‌های روی یال‌های $P = p(v_1, v_2)$ (مسیر) ادعا می‌کنیم که گره $v \in P$ وجود دارد به نحوی که یال‌های از هر دو گره v_1 و v_2 اوزان غیر نزولی به v دارند چرا که در غیر اینصورت یک یال کمینه محلی در مسیر پیدا می‌شد که باید حذف می‌شد.

اما حد اقل با یکی از v_1 و v_2 متفاوت است فرض کنید $v \neq v_1$ حال گره همسایه v_1 (مثلاً v_3) یال (v_1, v_3) را علامت زده است چون اوزان یال‌ها به سمت v فقط زیاد می‌شوند. پس توسط v_1 علامت زده شده است. پس (v_1, v_3) همان $c(v_1)$ طبق تعریف است. در عین حال (v_1, v_3) یک حد بالایی برای برش کمینه مابین v_1 و v_2 است پس لم ثابت می‌شود.

$$c(v_1, v_2) \geq \min(c(v_1), c(v_2))$$

□

الگوریتم های تقریبی برای پیدا کردن اجتماعات

در این فصل ما به روش Charikar در [۱۷] اشاره می کنیم و این روش را روی گراف وب با خصوصیات خاص خودش (رجوع کنید به فصل ۲) بررسی می کنیم. تا اینجا روشن است که

(۱) تعریف مشخصی برای اجتماعات وب وجود ندارد و هر تعریف ارائه شده جنبه هایی از اجتماعات را تاکید و جنبه هایی را نادید می گیرد.

(۲) لزومی به حل مساله اجتماعات بصورت دقیق وجود ندارد و راه حل های تقریبی با ضریب تقریب مناسب و مرتبه پایین به راه حل های دقیق مرتبه بالا ارجح می باشند.

در این بخش به بیان راه حل های تقریبی برای هر دو نوع گراف (بی جهت و جهت دار) می پردازیم و یک پیاده سازی کارا را برای راه حلمان شرح می دهیم.

۱-۵ گراف بدون جهت

تعریف ۱۸. اگر $G(V, E)$ گراف بدون جهت و $S \subseteq V$ باشد. مجموعه یالی القایی S یا $E(S)$ بصورت زیر تعریف می شود:

$$E(S) = \{ij \in E : i \in S, j \in S\}$$

تعریف ۱۹. اگر $S \subseteq V$ باشد. چگالی $f(S)$ این زیرمجموعه S از V بصورت زیر تعریف می شود:

$$f(S) = \frac{|E(S)|}{|S|}$$

در نهایت چگالی یک گراف بدون جهت $G(V, E)$ تعریف می شود:

$$f(G) = \max_{S \subseteq V} \{f(S)\}$$

لم ۲۰. $f(S)$ میانگین درجات زیر گراف القایی توسط S و $2f(G)$ بیشینه میانگین درجات روی تمام زیرگراف های القایی در G می باشد.

اثبات. جمع تمام درجات روی رئوس برابر $2E$ می باشد. پس حکم ثابت می شود. \square
مساله محاسبه $f(G)$ به مساله شلوغ ترین زیرگراف مشهور است که توسط ترفندهای الگوریتم شماره قابل حل است. یک مساله مشابه مساله شلوغ ترین k -زیرگراف است که روی تعداد گره ها در بیشینه $f(G)$ قید k می گذارد. این دو مساله ی بسیار جالب برای ارائه راه حل های تقریبی می باشد و نسبتاً باز هستند.

۵-۱-۱ محاسبه $f(G)$

محاسبه $f(G)$ بوسیله ترفندهای برنامه نویسی خطی (Linear programming) قابل انجام است. Charikar در [۱۷] علاوه بر این محاسبه یک الگوریتم 2-approximation برای این کار ارائه کرده

است. این کار بوسیله حذف متناوب گره با کمترین درجه صورت می گیرد. در ادامه مشابه همین بحث را برای گراف های جهت دار انجام می دهیم.
LP زیر را در نظر بگیرد.

$$\max \sum_{ij} x_{ij}$$

$$\forall ij \in E : x_{ij} \leq y_i$$

$$\forall ij \in E : x_{ij} \leq y_j$$

$$\sum_i y_i \leq 1$$

$$x_{ij}, y_i \geq 0$$

لم ۲۱. برای هر زیرمجموعه مقدار LP فوق حداقل برابر $f(S)$ است.

اثبات. اگر ما یک جواب امکانپذیر به این LP با مقدار $f(S)$ بدهیم حکم ثابت است. اگر $x = \frac{1}{|S|}$ باشد برای هر $i \in S$ قرار می دهیم $\hat{y}_i = x$ و برای هر $ij \in E(S)$ قرار می دهیم $\hat{x}_{ij} = x$ و سایر متغیرها را برابر صفر قرار می دهیم حال $\sum_i \hat{y}_i = |S|.x = 1$ پس (\hat{x}, \hat{y}) شرایط مساله را ارضاء می کنند. جواب این LP برابر می باشد با

$$|E(S)|.x = \frac{|E(S)|}{|S|} = f(S)$$

□ لم زیر نشان می دهد که همواره $S \in V$ وجود دارد که چگالی آن بیشتر از هر جواب امکانپذیر برای LP مذکور باشد.

لم ۲۲. با داشتن هر جواب امکانپذیر از LP فوق با مقدار v می توان $S \in V$ ساخت بشرطی که

$$f(S) \geq v$$

اثبات با پارامتریزه کردن مجموعه S و برهان خلف بدست می آید که در [۹] بطور کامل آمده است.

قضیه ۲۳. اگر مقدار جواب بهینه LP فوق را با $OPT(LP)$ نشان دهیم.

$$\max_{S \in V} \{f(S)\} = OPT(LP)$$

همچنین زیرمجموعه S برای جواب بهینه از روی راه حل بهینه LP قابل استخراج است.

۲-۵ الگوریتم تقریبی با ضریب ۲ برای محاسبه $f(G)$

هدف ما این است که در گراف G یک زیرگراف شلوغ بدست آوریم. طبیعتاً یک راه حل بدیهی این است که آن گره هایی از گراف که درجه کمتری دارند را بدور اندازیم تا شانس خود را برای بهبود چگالی بیازماییم. این ایده یک الگوریتم حریصانه را پیشنهاد می کند. مشابه چنین الگوریتمی برای میانگین درجه بزرگ برای تعداد گره های با عدد k توسط Asahiro و سایرین در [۱۸] تحلیل شده است.

برای این مساله ما از تحلیل Charikar در اثبات ضریب ۲ آن و پیاده سازی پیشنهادی او استفاده می کنیم. الگوریتم همواره یک زیرمجموعه $S \in V$ از گره های گراف را نگاه می دارد. در ابتدا $S \leftarrow V$ در هر دور الگوریتم i_{\min} یعنی گره با کمترین درجه تشخیص داده می شود الگوریتم را i_{\min} از مجموعه S حذف کرده و به دور بعدی می رود. شرط پایان الگوریتم زمانی است که S تهی می شود. برای تمام مجموعه های S که در طول اجرای الگوریتم ساخته شد آن مجموعه ای که بیشترین مقدار $f(S)$ را داشت بعنوان خروجی برگردانیده خواهد شد. در زیر این الگوریتم بصورت شبهه که آمده است.

```

ApproxCom( Graph G(V, E))
  initialize S = V;
  S_max = V;
  f_max = E/V;
  While(notEmpty(S))
    Let i_min be the vertex with the minimum degree in S;
    Remove i_min from S;
    Update S and E(S);
    f = E(S)/S;
    if (f > f_max)
      f_max = f;
      S_max = S;
  Return S_max;

```

قضیه ۲۴. الگوریتم فوق یک راه حل $\frac{2}{3}$ تقریب به مساله $f(G)$ می باشد.

اثبات. در حقیقت باید ثابت کنیم که

$$\max_{S \subseteq V} \{f(S)\} \leq \frac{2}{3}v$$

که v مجموعه ایست که ما در نهایت در الگوریتم بر می گردانیم. در وهله اول ما گراف بدون جهت را به صورت جهت دار نگاه می کنیم و $d(i)$ را درجه ورودی راس i در نظر می گیریم. پس در عمل به هر راس i یا z راس i یا z را نسبت داده ایم حال تعریف می کنیم.

$$d^{\max} = \max_i \{d(i)\}$$

حال اگر S_{\max} آن مجموعه S باشد که $f(S)$ را بیشینه می کند.

$$|E(S_{\max})| \leq |S_{\max}| \cdot d^{\max}$$

$$f(G) = \frac{|E(S_{\max})|}{|S_{\max}|}$$

حال انتساب یال ها به یکی از دو سرشان در حین الگوریتم بدین صورت انجام می شود: در ابتدا تمامی یال ها دست نخورده می باشند. وقتی i_{\min} در S حذف شد؛ تمامی یال هایش که به سایر رئوس S وصل هستند برچسب او را می خورند. ثابت حلقه بدین صورت است که کلیه یال ها بین رئوس در مجموعه کنونی S بدون برچسب می باشند در حالی که باقی یال ها تماما منتسب شده اند. پس در انتهای الگوریتم تمامی یال ها منتسب شده اند چون مجموعه S تهی شده است. حال ثابت می کنیم که

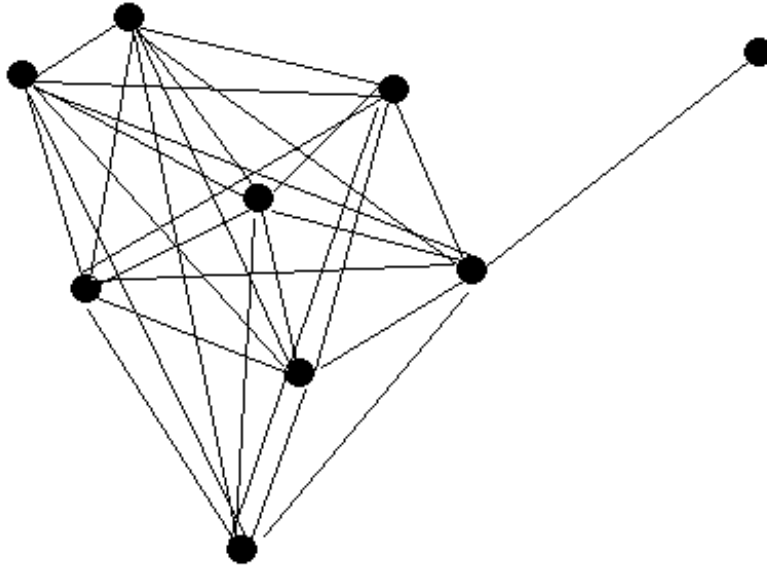
$$d^{\max} \leq 2v$$

که v بیشینه مقدار $f(S)$ است که توسط الگوریتم حریصانه بدست می آید. برای نشان دادن این مساله الگوریتم را در یک دور آن در نظر می گیریم. از آنجایی که i_{\min} طوری گزینش شده است که راس با کمترین درجه در S باشد پس درجه آن نمی تواند از $|E(S)|/|S|$ بیشتر باشد و می دانیم که این مقدار هم از $2v$ بیشتر نیست. همچنین در نظر داشته باشید که یک راس خاص تنها هنگامی یال بخود منتسب می کند که در حال حذف شدن باشد. پس الگوریتم ما یک تقریب ضریب ۲ برای $f(G)$ می باشد. \square

۳-۵ پیاده سازی الگوریتم

براحتی می توان یک پیاده سازی $O(n^2)$ برای الگوریتم فوق که n تعداد گره های گراف و m تعداد یال های آن باشد پیشنهاد کرد. بدین منوال که ما درجه رئوس را در زیر گراف القایی S نگاه می داریم. در هر دور از الگوریتم راس با کمترین درجه تشخیص و حذف می شود و درجه سایر رئوس باقیمانده بروز می شوند. هر دو این اعمال در $O(n)$ قابل انجام هستند پس الگوریتم در $O(n^2)$ به اتمام می رسد.

در اینجا ما پیشنهاد Charikar را برای اجرای این الگوریتم در $O(n)$ ارائه می کنیم. برای این کار n لیست از گره ها تهیه می کنیم که لیست i ام حاوی عناصر از درجه i می باشد. در هر دور از



گراف ناهنجار 5-1:

الگوریتم عنصر با کمترین درجه حذف شده و درجات عناصر همسایه اش بروز می شوند. این کار باعث انتقال این عنصر به لیست های دیگر می شود. توجه کنید که اگر حذف این راس باعث شدن i_{\min} گرهی در دور بعدی شده باشد ما کافی است که فقط با فاصله یک از درجه i_{\min} همین دور بدنبال i_{\min} دور بعد بگردیم به بیان دیگر اگر d درجه راس i_{\min} می باشد در دور بعد ابتدا در لیست d و $d+1$ و $d-1$ را نگاه می کنیم. البته می تواند این سه لیست خالی باشد ولی هنوز کار ما به اتمام نرسیده باشد به عنوان مثال گراف را 5-1 نگاه کنید.

با حذف i_{\min} در دور اول هر کدام از ۸ راس باقیمانده با درجه ی ۷ می تواند انتخاب شود. پس لزوما جستجو در $d-1$, d و $d+1$ کافی نیست. بلکه باید درجه اولین لیست غیرتهی را ثبت کنیم که اگر $d-1$ و d و $d+1$ باشد با جستجو حاصل می شود. در نهایت اگر مرتبه اینگونه جستجوها را ثابت بگیریم جواب در $O(m+n)$ حاضر می شود.

۴-۵ انتخاب d^{\max} بعنوان حد بالایی

در ارائه اثبات ضریب تقریب ایده حضور نابدیهی است در حالی که خود را در دوگان LP مساله ظاهر می کند. LP زیر دوگان LP مساله فوق است:

$$\min \gamma$$

$$\forall ij \in E : \alpha_{ij} + \beta_{ij} \geq 1$$

$$\forall i : \gamma \geq \sum_j \alpha_{ij} + \beta_{ji}$$

$$\alpha_{ij}, \gamma \geq 0$$

۵-۵ مساله ی معادل برای گراف جهت دار

ابتدا به بیان تعاریف معادل می پردازیم. اگر $G(V, E)$ یک گراف جهت دار و $S, T \subseteq V$ ، $E(S, T)$ مجموعه ای از یال هاست که از S به T می رود به بیان دیگر

$$E(S, T) = \{ij \in E : i \in S, j \in T\}.$$

پس بدین ترتیب چگالی جفت مجموعه S, T که آن را با $d(S, T)$ نمایش می دهیم بصورت زیر تعریف می شود.

$$d(S, T) = \frac{|E(S, T)|}{\sqrt{|S||T|}}$$

و مشابه گراف های بدون جهت چگالی $d(G)$ را برای یک گراف جهت دار $G(V, E)$ را تعریف می کنیم.

$$d(G) = \max_{S, T \subseteq V} \{d(S, T)\}$$

LP که برای مساله جدید تعریف می شود بصورت زیر است.

$$\max \sum_{ij} x_{ij}$$

$$\forall ij \in E : x_{ij} \leq s_i$$

$$\forall ij \in E : x_{ij} \leq t_j$$

$$\sum_i s_i \leq \sqrt{c}$$

$$\sum_j t_j \leq \frac{1}{\sqrt{c}}$$

$$x_{ij}, s_i, t_j \geq 0$$

برای گراف جهت دار هم خصوصیات مشابه قبیل قابل اثبات است. در اینجا ما این خصوصیات را بدون بیان اثبات می آوریم.

(۱) اگر $S, T \subseteq V$ و $c = \frac{|S|}{|T|}$ باشند. مقدار بهینه $LP(c)$ تعریف شده در بالا حداقل برابر $d(S, T)$ می باشد.

(۲) با داشتن یک جواب امکان پذیر از $LP(c)$ با مقدار v ما می توانیم زیرمجموعه های S, T را از V طوری بیابیم که $d(S, T) \geq v$

$$\max_{S, T \subseteq V} \{d(S, T)\} = \max_c \{OPT(LP(c))\} \quad (۳)$$

(۴) مجموعه S, T که بیشترین مقدار $d(S, T)$ را به ما می دهد از روی حل بهینه $LP(c)$ قابل محاسبه است.

۶-۵ الگوریتم حریمانه برای محاسبه $d(G)$

مشابه الگوریتمی که برای $f(G)$ توضیح داده شد برای یافتن $d(G)$ بکار می رود با این تفاوت که درجه ی یک گره در S یال هایی است که از آن گره به مجموعه T متصل است و درجه ی یک گره در T یال هایی است که از گره های موجود در S به آن اشاره شده است. در نهایت راس حذفی ما باید با یک استراتژی از بین گره های S یا T انتخاب شود. جفت زیرمجموعه S, T که مقدار $d(S, T)$ را بیشینه می کند بعنوان خروجی الگوریتم برگردانده می شود. Charikar نشان می دهد که الگوریتم فوق یک راه حل تقریبی با ضریب ۲ برای محاسبه $d(G)$ می باشد. البته یک حدس خوب برای مقدار c برای کارکرد مناسب الگوریتم لازم است. برای یک c خاص الگوریتم قابل اجرا در $O(m + n)$ می باشد. جستجو از مرتبه $O(\frac{\log n}{\epsilon})$ برای C لازم است که یک جواب $\epsilon + 2$ تقریب بدست آید.

نتایج تجربی روی الگوریتم تقریبی

Charikar

در ارتباط با سایر روش های شناسایی اجتماعات وب نتایج تجربی به تفصیل در کارهای اخیر یافت می شود. به عنوان نمونه به [۱۳] مراجعه کنید. با این وجود، روش تقریبی Charikar به منظور شناسایی اجتماعات وب تا کنون آزمایش نشده است.

الگوریتم تقریبی هر چند تضمین مناسبی از جهت زمان اجرا در اختیار ما می نهد، تنها هنگامی بکار شناسایی اجتماعات می آید که متریک های یک اجتماع را برآورده سازد.

در این بخش نتایجی که از اجرای الگوریتم تقریبی Charikar در فصل ۵ روی مجموعه داده واقعی وب بدست آمده را گزارش و بررسی می کنیم. در ابتدا تغییرات اعمالی روی الگوریتم تقریبی را شرح داده، سپس به تحلیل نتایج می پردازیم.

۱-۶ الگوریتم تقریبی استخراج اجتماعات وب

۱-۱-۶ اگر کل گراف در حافظه اصلی برنامه گنجانیده نشود

الگوریتم توضیح داده شده در فصل پیشین با فرض گنجانیده شدن کل گراف ورودی در حافظه اصلی برنامه می‌باشد. همانطور که نشان دادیم، چنین فرضی لزوماً در مورد گراف غول‌آسای وب صحیح نیست. پس در عمل الگوریتم زیر را برای شناسایی اجتماعات گراف وب پیشنهاد می‌کنیم.

```

1 Approx-Huge-Web (set: S, integer: k)
2   while number of iteration is less than desired do
3       set G to a crawl from S of depth k
4       set X = Approx-com ( G );
5       Score all the vertices in X due to their degree
6       Add non-seed high-scores in X to S
7   return X

```

رویه فوق بدین صورت است که در ابتدا یک مجموعه دانه که قابل بارگذاری در حافظه اصلی است را در نظر می‌گیریم. تا عمق k این مجموعه را خزش می‌کنیم و گراف ورودی ما به رویه Approx-com در فصل ۵ چنین ساخته می‌شود.

حال اگر X خروجی رویه Approx-com باشد، رئوس با درجه بالا که در S نبودند به S در مرحله بعد اضافه می‌شوند. در پیوست کد کامل به زبان جاوا آمده است.

۲-۱-۶ اگر بخواهیم همه اجتماعات را شناسایی کنیم

الگوریتم Approx-com توضیح داده شده در فصل قبل، تنها شلوغ‌ترین زیرگراف را برمی‌گرداند. برای شناسایی کلیه اجتماعات می‌توان الگوریتم زیر را ارائه داد.

Approx-all-communities (Graph (G,V), int n)

```

while i<n
    compute the ith community from Approx-com (G)
    let S_i be the output of Approx-com (G)
    remove S_i from G
return S_1, S_2, ... , S_n

```

در الگوریتم فوق فرض بر این است که کل گراف در حافظه اصلی نگهداری می‌شود.

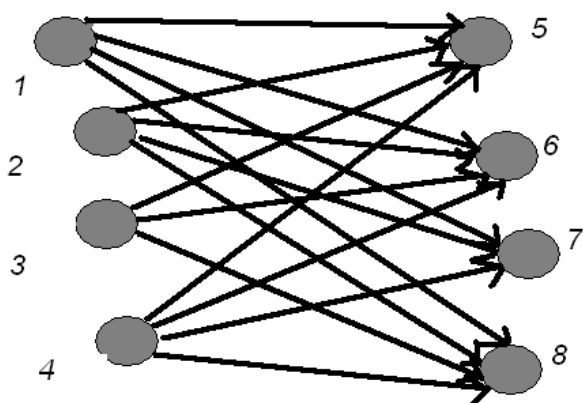
۲-۶ داده های واقعی

برای محک الگوریتم تقریبی Charikar از یک مجموعه ۶۰۰۰۰ راسی که خود یکی از ۱۰ مجموعه خزش و جمع‌آوری شده از صفحات وب، در دانشگاه شریف است. این مجموعه در سال ۲۰۰۴ توسط تیم طراحی موتور جستجو در آزمایشگاه دکتر قدسی، دانشکده مهندسی کامپیوتر خزش شده است.

برای اجرای الگوریتم، از یک مجموعه ۸۰۰ گرهی به منظور S یا مجموعه دانه استفاده کردیم. k را در الگوریتم برابر ۱ قرار دادیم. گراف شکل ۶-۱ به عنوان خروجی الگوریتم بدست آمد.

صفحات متناظر با هر گره در زیر آمده است.

- 1 <http://editors.dmoz.org/editors/editcat.cgi?cat=Computers/Security>
- 2 http://p2p.wrox.com/default.asp?CAT_ID=27
- 3 <http://s.teoma.com/search?q=books>
- 4 <http://www.altavista.com/cgi-bin/query?q=programming\%20languages>
- 5 <http://www.pcisystems.com>
- 6 <http://www.phytec.com/>



گراف خروجی الگوریتم تقریبی روی مجموعه دانه‌ی ۸۰۰ تایی 6-1:

7 <http://www.pond.ie>

8 <http://www.square1industries.com>

الگوریتم در دور دوم با حذف گره‌های دور اول به کار خود ادامه می‌دهد. اجتماع تشخیص داده شده در دوم در گراف ۶-۲ نمایش داده شده است. صفحات متناظر با هر گره در زیر آمده است.

1 : dmoz.org/Games/Video_Games/

[Console_Platforms/Sony/PlayStation/index.html](http://dmoz.org/Games/Video_Games/Console_Platforms/Sony/PlayStation/index.html)

2 : www.vidgames.com

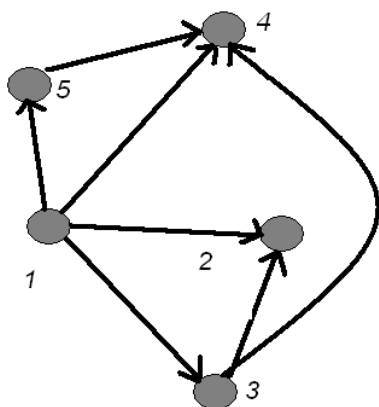
3 : www.playstationultra.com

4 : directory.google.com/Top/Games/

[Video_Games/Console_Platforms/Sony/PlayStation/](http://directory.google.com/Top/Games/Video_Games/Console_Platforms/Sony/PlayStation/)

5 : dmoz.org/cgi-bin/apply.cgi?where=

[Games2FVideo_Games2FConsole_Platforms2FSony2FPlayStation](http://dmoz.org/cgi-bin/apply.cgi?where=Games2FVideo_Games2FConsole_Platforms2FSony2FPlayStation)



گراف خروجی الگوریتم تقریبی روی مجموعه دانه‌ی ۸۰۰ تایید در دور دوم: 6-2

۳-۶ تحلیل نتایج

نتایج بدست آمده توسط الگوریتم تقریبی، انتظارات ما را در شناسایی درست اجتماعات برآورده می‌سازد. میزان‌های bibliometric و co-citation ارائه شده در الگوریتم HITS را بر روی مجموعه جواب محاسبه می‌کنیم. برای زیرگراف شکل ۶-۱ متریک‌ها بصورت زیر می‌باشد.

$$Auth - score(\Delta) = \sum_{i=1,2,3,4} Hub - score(i)$$

$$Auth - score(\Gamma) = \sum_{i=1,2,3,4} Hub - score(i)$$

$$Auth - score(\Upsilon) = \sum_{i=1,2,4} Hub - score(i)$$

$$Auth - score(\Lambda) = \sum_{i=1,2,3,4} Hub - score(i)$$

$$Hub - score(\uparrow) = \sum_{i=5,6,7,8} Auth - score(i)$$

$$Hub - score(۲) = \sum_{i=۵,۶,۷,۸} Auth - score(i)$$

$$Hub - score(۳) = \sum_{i=۵,۶,۸} Auth - score(i)$$

$$Hub - score(۴) = \sum_{i=۵,۶,۷,۸} Auth - score(i)$$

در دور دوم هم یک اجتماع بامعنی استخراج شده است.

چند نکته قابل توجه در شکل های ۱-۶ و ۲-۶ در زیر آمده است.

- نسبت تعداد یال به راس در گراف ۱-۶ به مقدار $\frac{۱۵}{۸}$ یا ۱.۸۷۵ می باشد.
- در گراف ۱-۶ یک هسته ی دوبخشی چهار در سه وجود دارد.
- گراف ۱-۶ با دریافت یک یال دیگر به گراف کامل دوبخشی $K_{۴,۴}$ تحویل می شود.
- نسبت تعداد یال به راس در گراف ۲-۶ تا مقدار $\frac{۱۶}{۷}$ یا ۱.۲ کاهش یافته است.
- با توجه به URL های ذکر شده، هر دو جواب، دو اجتماع صحیح است که گرد دو موضوع معنایی پدیدار شدند.
- اعمال متریک های دوبخشی در مورد ۲-۶ هم نتایج قابل قبولی می دهد.

در دوره های بعدی الگوریتم، باز هم $f(G)$ کاهش پیدا می کند. با وجودی که الگوریتم تقریبی اساساً بر پایه اجتماعات متراکم عمل می کند، در این مثال مشاهده کردیم که یک جواب خوب به مساله اجتماعات دوبخشی نیز می تواند باشد که مرتبه زمانی بسیار مطلوبی نسبت به روش های پیشین دارد.

نتیجه گیری

در این پایان نامه به بررسی روش های مختلف تشخیص اجتماعات موجود در وب پرداخته شده است. در ابتدا اهمیت این مساله، جنبه های مختلف آن توضیح داده شده و چند مساله مشابه معرفی شده اند. در ادامه خصوصیات اصلی گراف وب را بر شمرده ایم و یک مدل ریاضی برای شبیه سازی این گراف ارائه کردیم.

در ادامه به تشریح الگوریتم های مختلف تشخیص اجتماعات وب پرداختیم که چندی از نتایج باختصار در زیر آمده است.

- متریک ها و میزان هایی برای تعیین میزان شباهت دو صفحه ای که به یکدیگر اتصال مستقیم ندارد تعریف می شود.
- هسته های دو بخشی براساس دید هاب و مرجع بودن در یک گراف اساس تعریف اجتماعات وب در روش HITS است.
- هر دو روش HITS و PageRank تفسیرهای ماتریسی و قدم زدن تصادفی دارند و به یکدیگر قابل تبدیل هستند.
- روش پیشینه شاره دید متفاوتی نسبت به HITS و PageRank از اجتماعات وب دارد.
- روش درخت برش کمینه یک روش ساخت یافته براساس الگوریتم شاره است.
- الگوریتم تقریبی نتایج قابل قبولی با پیچیدگی خطی برای پیدا کردن اجتماعات می دهند.

مسائل فراوانی را هر کدام از روش های تشخیص اجتماعات پیش روی ما باز می گذارند. در روش تقریبی Charikar یک زمینه تحقیقات فراتر بکارگیری خصوصیات گراف وب معرفی شده در بخش نخست است برای کاهش ضریب تقریب تئوریک و یا بهبود الگوریتم در عمل. حل مساله با تعمیم تعریف اجتماعات در روش شماره می تواند یک چالش جدی در این عرصه باشد. همچنین استفاده از اطلاعات متنی علاوه بر ساختار پیوندی گراف وب یک پیشنهاد عملی خوب است.

فصل ۸

واژه‌نامه

Communities	اجتماعات
feasible	امکان پذیر
contraction	انقباض
query	پرسش
approximation	تقریب
singleton	تک عضوی
power law distribution	توزیع توانی
strongly connected component	جزء قویا همبند
weakly connected component	جزء همبند ضعیف
sink	چاهک
source	چشمه
greedy	حریصانه
crawl	خزش
domain	دامنه
cluster	دسته
clustering	دسته بندی

seed	دانه
link structure	ساختار پیوندی
flow	شاره
Webpage	صفحه وب
minimal	غیر قابل کاهش
random walker	قدم‌زن تصادفی
search term	کلمه کاوشی
node	گره
dense	متراکم
authority	مرجع
average connected distance	میانگین فاصله اتصالی
host	میزبان
bipartite core	هسته دوبخشی
edge	یال

کتاب نامه

- [1] K. Bharat and M.R. Henzinger. Improved algorithms for topic distillation in a hyperlinked environment. *Proc. ACM SIGIR, 1998.*
- [2] R. Albert, H. Jeong, and A. Barabasi. Diameter of the World Wide Web .Nature 401:130-131, 1999.
- [3] S. Brin and L. Page. The anatomy of a large-scale hyper textual Web search engine. *Proc. 7th WWW Conf. 1998*
- [4] N. Gilbert. A simulation of the structure of academic science. *Sociological Research Online, 2(2), 1997.*
- [5] K. Bharat, B. Chang, M. Henzinger, M. Ruhl Who links to whom: Mining Linkage between Web Sites
- [6] R. Kumar, P. Raghavan, S. Rajagopalan, D. Sivakumar, A. Tomkins, E. Upfal. The Web as a graph.
- [7] A. Mendelzon, G. Mihaila, and T. Milo. Querying the World Wide Web. *J. of Digital Libraries, 1(1):68-88, 1997.*
- [8] S. R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Trawling emerging cyber-communities automatically. *Proc. 8th WWW Conf, 1999.*

- [9] S. R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tomkins. Extracting large-scale knowledge bases from the web. *Proc. VLDB, 1999*.
- [10] J. M. Kleinberg. *Authoritative sources in a hyperlinked environment*. In Proceedings of the ninth annual ACM-SIAM Symposium on Discrete Algorithms, pages 668-677, 1998.
- [11] A. V. Goldberg, E. Tardos, and R. E. Tarjan. Network Flow algorithms. In B. Korte, L. Lovasz, H. J. Promel, and A. Schrijver, editors, *Paths, flows, and VLSI-layout, volume 9 of Algorithms and Combinatorics*, , pages 101-164. Springer-Verlag, 1990.
- [12] R. E. Gomory and T. C. Hu. Multi-terminal network flows. *J. SIAM*, 9:551-570, 1961.
- [13] K. Tsioutsoulis. *Maximum Flow Techniques for Network Clustering*. PhD thesis, Princeton University, Princeton, NJ, June 2002.
- [14] G. W. Flake, R. E. Tarjan, and K. Tsioutsoulis. *Graph clustering and minimum cut trees*. Internet Mathematics, 2003.
- [15] A. Benczur. *Counterexamples for directed and node capacitated cut-trees*. *SIAM J. Comput.* 24:3 (1995), pp. 505-510.
- [16] Z. Wu and R. Leahy. An optional graph theoretic approach to data clustering: Theory and its application to image segmentation. *IEEE Transaction on Pattern Analysis and Machine Intelligence*. November 1993.
- [17] M. Charikar. *Greedy Approximation Algorithms for Finding Dense Components in Graphs*, in *Proceedings of APPROX, 2000*.
- [18] Y. Asahiro, K. Iwama, H. Tamaki and T. Tokuyama. Greedily Finding a Dense Subgraph. *Journal of Algorithms*, 34(2):203-221, 2000.

Web Communities and Algorithms

Abstract

Communities are those sets of entities in Web that have a strong relationship to each others. They appear in the shape of Web pages or Host pages. It is necessary to examine the link structure of the induced graph of Web in order to infer the communities. Recognizing Web communities provide us with a precious information about what pages are about.

Given a graph $G(V, E)$ with V representing the Web pages and E representing the links among them, the problem is to extract some subsets from G , which are either in the form of hubs and authorities, or dense in edges.

In this B.Sc thesis, we probe the web characterisitcs and models. We explain the traditional methods to identify Web communities. We will see the minimum cut tree method. Then we examine the Charikar's procedure and then we conclude.

Keywords: *Communities, hub, authority, flow, cut*

Web Communities and Algorithms

by

Mahyar Salek

Submitted in Partial Fulfillment
of the Requirements
for the Degree of
Bachelor of Science
in
Computer Engineering (Software)

Under supervision of
Prof. Mohammad Ghodsi

August 2005

Computer Engineering Department
Sharif University of Technology
Tehran