



Sharif University of Technology

Department of Computer Engineering

M.Sc. Thesis

**Approximation Algorithms for Clustering Points
in the Distributed Model**

By:

Emad Aghajani

Supervisor:

Dr. Hamid Zarrabi-Zadeh

July 2016

Abstract

Clustering is one of the most well-known fundamental problems in computer science. In this thesis we have focused on a particular version of such problem, called capacitated k -center, and we analyze and survey them in the distributed model, when massive data is given. Our contribution in this research includes a new approximation algorithms with constant approximate factors for such problems in the distributed model, as well as improving the best available algorithm for capacitated k -center problem. Composable coresets as one of the most important techniques in distributed and streaming model is our primary tools in designing these algorithms. This technique gives the opportunity of solving the problem in smaller chunks of data, and giving the result by combining them.

Keywords: Approximation algorithm, Composable coresets, Clustering, Capacitated k -center, Distributed model



دانشگاه صنعتی شریف
دانشکده‌ی مهندسی کامپیوتر

پایان‌نامه‌ی کارشناسی ارشد
گرایش مهندسی نرم‌افزار

عنوان:

الگوریتم‌های تقریبی برای خوشه‌بندی نقاط در مدل توزیع شده

نگارش:

عماد آقاجانی

استاد راهنما:

دکتر حمید ضرابی‌زاده

مرداد ۱۳۹۵

صلاة الاضحية

به نام خدا
دانشگاه صنعتی شریف
دانشکده‌ی مهندسی کامپیوتر

پایان‌نامه‌ی کارشناسی ارشد

عنوان: الگوریتم‌های تقریبی برای خوشه‌بندی نقاط در مدل توزیع شده
نگارش: عماد آقاجانی

کمیته‌ی ممتحنین

امضاء: استاد راهنما: دکتر حمید ضرابی‌زاده

امضاء: استاد مشاور: دکتر محمد علی آبام

امضاء: استاد مدعو: دکتر علیرضا زارعی

تاریخ:

سپاس

نخستین سپاس و ستایش از آن خداوندی است که بنده کوچکش را در دریای بیکران اندیشه، قطره‌ای ساخت تا وسعت آن را از دریچه اندیشه‌های ناب آموزگارانی بزرگ به تماشا نشیند. لذا اکنون که در سایه‌سار بنده‌نوازی‌هایش پایان‌نامه حاضر به انجام رسیده است، بر خود لازم میدانم تا مراتب سپاس را از بزرگوارانی به‌جا آورم که اگر دست یاری‌گرشان نبود، هرگز این پایان‌نامه به انجام نمی‌رسید. ابتدا از استاد گران‌قدرم جناب آقای دکتر ضرابی‌زاده که زحمت راهنمایی این پایان‌نامه را بر عهده داشتند، کمال سپاس را دارم.

سپس از دوست عزیزم، مهندس کیان میرجلالی، که در به دست آوردن نتایج این مقاله همکاری داشتند، صمیمانه سپاس گزارم.

و سپاس ویژه به مهربان‌ترین همراهان زندگی‌ام، به پدر، مادر و برادر عزیزم تقدیم می‌کنم که حضورشان در فضای زندگی‌ام مصداق بی‌ریای سخاوت بوده است.

چکیده

مسائل خوشه‌بندی نقاط از مهم‌ترین و شناخته‌شده‌ترین مسائل پایه‌ای در علوم کامپیوتر است. در این پژوهش ما بر روی دسته‌ی خاصی از این مسائل که با نام « k -مرکز ظرفیت دار» شناخته می‌شوند تمرکز خواهیم کرد، و سعی می‌کنیم تا این دسته از مسائل را در مدل توزیع‌شده که با حجم زیادی از داده روبرو هستیم، بررسی نموده و الگوریتم‌های تقریبی جدیدی برای آن ارائه کنیم. در ارائه این الگوریتم‌ها از شیوه «مجموعه‌های هسته‌ی ترکیب‌پذیر» که در مدل‌های توزیع‌شده مورد توجه قرار گرفته است، کمک گرفته‌شده است. این مفهوم به ما کمک می‌کند تا به ازای داده‌های حجیم، مسئله را به بخش‌های کوچک شکسته و از اجتماع نتایج، برای ارائه یک جواب تقریبی استفاده کنیم. ارائه چند الگوریتم تقریبی با ضریب تقریب ثابت برای نسخه‌های مختلف مسئله k -مرکز ظرفیت دار در مدل توزیع‌شده و بهبود ضریب تقریب الگوریتم‌های فعلی، از نتایج این پژوهش به‌شمار می‌آید.

کلیدواژه‌ها: الگوریتم تقریبی، مجموعه‌های هسته‌ی ترکیب‌پذیر، خوشه‌بندی، k -مرکز ظرفیت دار، مدل توزیع‌شده

فهرست مطالب

۱۰	۱	مقدمه
۱۱	۱-۱	تعریف مسئله
۱۴	۲-۱	اهمیت موضوع
۱۴	۳-۱	ادبیات موضوع
۱۶	۴-۱	اهداف تحقیق
۱۷	۵-۱	ساختار پایان نامه
۱۸	۲	مفاهیم اولیه
۱۸	۱-۲	الگوریتم های تقریبی
۲۰	۲-۲	خوشه بندی
۲۳	۳-۲	مجموعه هسته
۲۵	۳	کارهای پیشین
۲۵	۱-۳	مسائل خوشه بندی
۲۸	۲-۳	مسئله k -مرکز در مدل توزیع شده
۳۱	۴	الگوریتم تقریبی جدید برای k -مرکز وزن دار
۳۲	۱-۴	تعاریف و مفاهیم اولیه

۳۴	۲-۴	الگوریتم پیشنهادی
۴۱	۳-۴	اثبات درستی الگوریتم
۴۶	۵	الگوریتم‌های تقریبی جدید برای مسئله k -مرکز ظرفیت‌دار در مدل توزیع شده
۴۶	۱-۵	تعاریف و مفاهیم اولیه
۴۸	۲-۵	الگوریتم تقریبی برای مسئله k -مرکز با ظرفیت‌های نرم
۴۸	۱-۲-۵	رابطه‌ی جواب مسئله k -مرکز و مسئله k -مرکز با ظرفیت‌های نرم
۴۹	۲-۲-۵	یک لم کلیدی
۵۱	۳-۲-۵	الگوریتم پیشنهادی
۵۳	۳-۵	الگوریتم بهبودیافته برای مسئله k -مرکز با ظرفیت‌های نرم
۵۵	۴-۵	الگوریتم تقریبی برای مسئله k -مرکز با ظرفیت‌های سخت
۵۸	۶	نتیجه‌گیری
۵۸	۱-۶	نتایج بدست آمده
۶۰	۲-۶	کارهای آینده

فهرست شکل‌ها

- ۱-۱ مقایسه شعاع بهینه مسئله k -مرکز برای حالت بی‌ظرفیت و ظرفیت‌دار ۱۳
- ۱-۲ مثالی از مجموعه هسته ۲۴
- ۱-۴ یک مثال برای معرفی نمادهای معرفی‌شده در این بخش ۳۳
- ۲-۴ مثالی از مجموعه پادشاه‌ها و امپراتوری ۳۷
- ۳-۴ مثالی از گراف دوبخشی G' ۳۸
- ۱-۵ رابطه هندسی بین سه نقطه $p \in \bar{S}$ و $\phi(p)$ و $\bar{\phi}(p)$ ۵۱
- ۲-۵ نحوه جایگزینی مرکز برای تبدیل جواب نسخه ظرفیت نرم به سخت ۵۷

فهرست جدول‌ها

- ۱-۲ تابع هدف تعدادی از مسائل خوشه‌بندی ۲۱
- ۱-۶ ضریب تقریب‌های بدست آمده برای نسخه‌های مختلف مسئله‌ی k -مرکز ظرفیت‌دار ۵۹
- ۲-۶ مقایسه الگوریتم پیشنهادی با نمونه موجود برای مسئله‌ی k -مرکز ظرفیت‌دار ۵۹

فصل ۱

مقدمه

در این فصل مقدمه‌ای از مسائل موردنظر در این پژوهش و اهداف و انگیزه‌های آن‌ها را ارائه خواهیم داد و به معرفی مسائل خوشه‌بندی و به‌طور خاص مسئله‌ی k -مرکز و نسخه‌های ظرفیت‌دار آن که هدف اصلی این پژوهش است می‌پردازیم. سپس به نحوی برخورد با این دسته از مسائل به ازای ورودی‌های بزرگ و پژوهش‌های انجام‌شده در مدل توزیع‌شده و به‌طور خاص با استفاده از مفهوم مجموعه‌های هسته، تمرکز خواهیم کرد.

مسئله‌ی خوشه‌بندی نقاط، یکی از مهم‌ترین و شناخته‌شده‌ترین دسته مسائل پایه‌ای است که کاربردهای آن به‌طور گسترده در صنعت مشاهده می‌شود. به‌عنوان نمونه می‌توان مجموعه‌ای از نقاط که فاصله بین نقاط میزان شباهت را نشان می‌دهد در نظر گرفت. در این صورت مسائل خوشه‌بندی به دنبال ارائه یک دسته‌بندی از نقاط خواهند بود طوری که اجسام مشابه در یک دسته قرار بگیرند و این موضوع می‌تواند در موتورهای جستجو مورد استفاده قرار بگیرد. مسائلی از قبیل k -مرکز و k -میانه، از جمله این مسائل خوشه‌بندی می‌باشند که در تعریف تابع هدف^۱ تفاوت‌های جزئی با یکدیگر دارند. همچنین در دسته‌ای از مسائل، عمل خوشه‌بندی با محدودیت‌های جدید همراه می‌شود. به‌عنوان نمونه، در نسخه‌های ظرفیت‌دار مسئله‌ی k -مرکز، برای اندازه هر خوشه یک محدودیت در نظر گرفته می‌شود.

از آنجا که اغلب مسائل خوشه‌بندی ان‌پی-سخت^۲ هستند [۱] و نمی‌توان در زمان چندجمله‌ای پاسخ

^۱ Objective Function

^۲ NP-hard

دقیقی برای آن‌ها به دست آورد، (مگر $P = NP$)، عموماً پژوهشگران به دنبال ارائه الگوریتم‌هایی برای حل تقریبی مسئله بوده‌اند.

از طرف دیگر، در سال‌های اخیر حجم داده‌های خامی که نیازمند پردازش می‌باشند به شکل گسترده‌ای افزایش یافته است. به همین دلیل استفاده از مدل‌های توزیع‌شده و چارچوب‌هایی همچون نگاشت-کاهش مورد استقبال قرار گرفته است. در این مدل‌ها سعی می‌شود تا مسئله توسط تعدادی ماشین مجزا و به صورت موازی حل شود. بین ماشین‌ها می‌تواند به صورت محدود تبادلات داده‌ای وجود داشته باشد. در این پژوهش قصد داریم تا چند روش تقریبی برای حل مسئله k -مرکز ظرفیت‌دار در فضای توزیع‌شده ارائه دهیم. این پژوهش بر دو نسخه از مسئله k -مرکز ظرفیت‌دار، نسخه با ظرفیت‌های نرم و نسخه‌ی با ظرفیت‌های سخت، تمرکز خواهد داشت. الگوریتم‌های ارائه‌شده در این پژوهش، یک ضریب تقریب بهبودیافته نسبت به الگوریتم‌های موجود ارائه خواهند داد.

۱-۱ تعریف مسئله

مسئله اصلی مورد بحث در این پژوهش مسئله k -مرکز است. هدف از مسئله k -مرکز انتخاب مجموعه‌ی P از k نقطه از میان مجموعه نقاط ورودی مسئله S است به طوری که فاصله دورترین نقطه از مجموعه نقاط P کمینه باشد^۳. هر یک از نقاط مجموعه P نماینده یک خوشه فرض می‌شوند و در اصطلاح «مرکز» یک خوشه نامیده می‌شوند. در بعضی فرمول‌بندی‌های خوشه‌بندی، خوشه‌بندی تنها با انتخاب مرکزها انجام می‌گیرد و نیازی به انتساب سایر نقاط به این مراکز نیست. چراکه اگر هر نقطه را به نزدیک‌ترین مرکز خود منتسب کنیم به یک خوشه‌بندی از نقاط دست پیدا خواهیم کرد؛ اما این حرف تنها تا زمانی درست است که محدودیتی بر اندازه هر خوشه وجود نداشته باشد و در نسخه‌های ظرفیت‌دار که مورد علاقه ما است، نیاز است تا نحوه انتساب نیز به طور دقیق مشخص گردد.

با این فرض که مسئله در فضای متری^۴ مورد بحث است، می‌توان ورودی مسئله را یک گراف تصور نمود. در این صورت، مسئله k -مرکز به دنبال انتخاب k رأس از گراف و انتساب سایر رئوس به آن‌ها است به طوری که بیشترین فاصله از یک رأس تا مرکز منتسب شده به آن که آن را شعاع خوشه می‌نامیم،

^۳ فاصله یک نقطه مانند p از یک مجموعه نقطه S برابر فاصله نقطه p تا نزدیک‌ترین نقطه از مجموعه S است

^۴ Metric

کمینه باشد. این مسئله به صورت رسمی به صورت زیر تعریف می شود:

تعریف ۱-۱ (مسئله k -مرکز) به ازای گراف وزن دار $G = (V, E)$ ، زیرمجموعه‌ی $S \subseteq V$ با اندازه حداکثر k پیدا کنید طوری که عبارت زیر کمینه شود.

$$\max_{u \in V} \min_{v \in S} d(u, v)$$

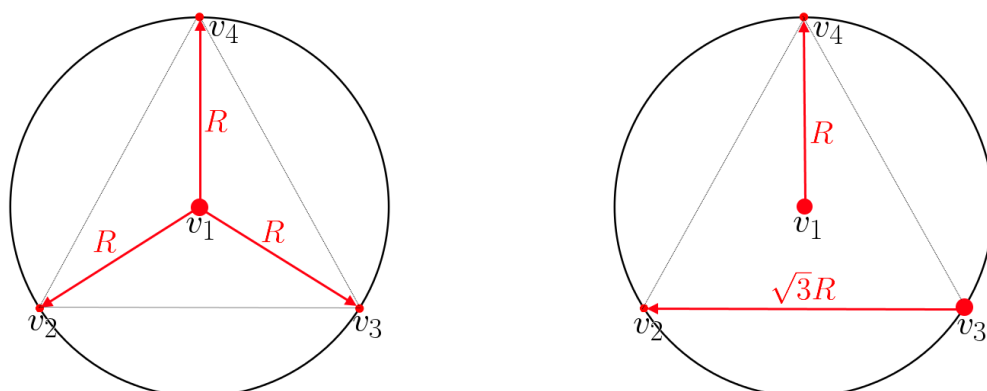
مسئله‌ی k -مرکز به شکل‌های مختلف قابل تعمیم است. در دسته‌ای از مسائل به ازای هر نقطه از ورودی به عنوان یک مرکز، یک هزینه بازگشایی برای آن نقطه در نظر گرفته می شود که به مسئله «مکان‌یابی تسهیلات^۵» مشهور است. در دسته‌ای دیگر از پژوهش‌ها، هدف کم کردن تعداد مراکز با یک شعاع ثابت است. دسته‌ای از پژوهش‌ها نیز به اعمال ظرفیت بر روی اندازه هر خوشه می پردازند و به مسئله « k -مرکز ظرفیت دار» مشهور است که ما در این پژوهش به طور خاص بر روی آن تمرکز خواهیم داشت. همچنین در نسخه‌ی خاصی از مسئله‌ی k -مرکز ظرفیت دار، به ازای هر نقطه یک وزن قائل می شویم. در این صورت پوشش یک نقطه، به اندازه وزن آن نقطه، از ظرفیت مرکز استفاده می کند. این نسخه از مسئله به «نسخه وزن دار مسئله‌ی k -مرکز ظرفیت دار» مشهور است. این نسخه از مسئله به عنوان ابزاری برای ارائه الگوریتم نهایی مورد استفاده قرار خواهد گرفت و در بخش ۴-۲ به تفصیل به آن پرداخته خواهد شد.

در مسئله k -مرکز ظرفیت دار^۶ محدودیت حداکثر تعداد نقاط مجاز در هر خوشه (حداکثر نقاط قابل انتساب به یک مرکز) به مسئله اصلی اضافه می شود. این نسخه با انگیزه‌ی هرچه کاربردی تر شدن مسائل مورد توجه قرار گرفته است. برای نمونه فرض کنید که ما به دنبال قرار دادن دو آنتن ($k = 2$) بر روی مجموعه رئوس $V = \{v_1, v_2, v_3, v_4\}$ هستیم، به طوری که شعاع لازم برای پوشش تمام نقاط کمینه باشد (همانند مسئله‌ی k -مرکز بیان شده).

پاسخ بهینه برای این مسئله (شعاع برابر با R) در شکل ۱-۱-۱ سمت چپ نمایش داده شده است؛ اما در دنیای واقعی یک آنتن ظرفیت محدودی در سرویس دهی دارد. اگر ظرفیت تمام آنتن‌ها را، به صورت یکسان، یک واحد فرض کنیم ($L = 1$)، در این صورت جواب بهینه مسئله متفاوت خواهد بود (شعاع برابر با $\sqrt{3}R$) که در شکل ۱-۱-۱ سمت راست نمایش داده شده است.

^۵ Facility Location

^۶ Capacitated k -Center



شکل ۱-۱: مقایسه شعاع بهینه مسئله k -مرکز برای حالت بی ظرفیت و ظرفیت دار

مسئله k -مرکز ظرفیت دار به صورت رسمی به صورت زیر تعریف می شود:

تعریف ۱-۲ (مسئله k -مرکز ظرفیت دار) به ازای گراف وزن دار $G = (V, E)$ ، زیرمجموعه $S \subseteq V$ با اندازه حداکثر k پیدا کنید طوری که حداکثر L رأس به هر مرکز از S منتسب شده باشد و عبارت زیر نیز کمینه شود.

$$\max_{u \in V} \min_{v \in S} d(u, v)$$

تاکنون نسخه های ظرفیت دار با تعاریف متعددی در پژوهش های مختلف مورد مطالعه قرار گرفته اند. خولر و سوسمان [۲] دو نسخه ظرفیت دار پایه ای از این مسئله با نام های « k -مرکز با ظرفیت های نرم^۷» و « k -مرکز با ظرفیت های سخت^۸» را مورد بررسی قرار داده اند.

در تعریف مسئله k -مرکز با ظرفیت های سخت، ظرفیت هر خوشه حداکثر L در نظر گرفته می شود. در مسئله k -مرکز با ظرفیت های نرم در کنار محدودیت حداکثر L تایی بر روی هر خوشه، این اجازه داده می شود تا یک نقطه بیش از یک بار به عنوان مرکز استفاده شود. در این صورت ما می توانیم در هر نوبت از الگوریتم، بدون در نظر گرفتن انتخاب های قبلی، هر یک از نقاط ورودی را به عنوان مرکز انتخاب کنیم.

^۷ k -Center with soft capacities / Capacitated multi- k -Center

^۸ k -Center with hard capacities / L-Balanced k -Center

۲-۱ اهمیت موضوع

حجم گسترده مطالعات انجام شده بر روی دسته مسائل خوشه‌بندی از یک سو، و میزان توجه پژوهشگران به حالت‌های خاص و نسخه‌های گوناگون مسئله از سوی دیگر، می‌تواند اهمیت این طیف از مسائل را در دنیای علوم کامپیوتر به خوبی نشان دهد. در رابطه با مسئله k -مرکز به‌طور خاص، کاربردهای عملی گوناگونی وجود دارد که از جمله آن‌ها می‌توان به موضوع مکان‌یابی تسهیلات، دسته‌بندی داده‌ها و موتورهای جستجو اشاره نمود. همچنین گفتنی است که دسته‌ای از پژوهش‌ها به‌طور خاص به بررسی کاربردهای مسئله k -مرکز پرداخته‌اند [۳].

با توجه به آنکه در بعضی از کاربردها مانند مثالی که در بخش قبل ارائه شد، وجود تعادل در خوشه‌بندی ضروری است، نسخه‌های ظرفیت‌دار از مسئله k -مرکز مطرح شده‌اند. همچنین تعداد پژوهش‌های صورت گرفته در یک دهه اخیر بر روی این دسته از مسائل نشان‌دهنده توجه ویژه پژوهشگران و صنعتگران به این مسائل است. همچنین از آنجاکه در عمل یکی از بزرگ‌ترین کاربردهای مسائل خوشه‌بندی برای تحلیل مجموعه داده‌های حجیم^۹ است و امکان پردازش چنین حجم از داده‌هایی به صورت ترتیبی و متمرکز وجود ندارد، پژوهشگران به دنبال روش‌های حل توزیع شده و یا تقریبی مسئله می‌باشند.

۳-۱ ادبیات موضوع

با توجه به آنکه مسئله اصلی مورد بحث در این پژوهش ارائه یک الگوریتم تقریبی برای حل مسئله k -مرکز ظرفیت‌دار در مدل توزیع شده است، ما با دو طیف گسترده از پژوهش‌ها در ارتباط خواهیم بود. دسته اول پژوهش‌هایی هستند که به مسائل خوشه‌بندی و به‌طور خاص به مسئله k -مرکز و نسخه‌های ظرفیت‌دار آن توجه دارند و دسته دوم نیز مطالعات مربوط به مدل‌ها و شیوه‌های توزیع‌شدگی از قبیل مدل نگاشت-کاهش است. همچنین آن دسته از پژوهش‌ها که به ترکیب این دو پرداخته‌اند نیز مورد علاقه ما است.

پیش‌تر ثابت شده است که دسته مسائل خوشه‌بندی غالباً NP -سخت^{۱۰} هستند [۱] و نمی‌توان در

^۹ Massive Data

^{۱۰} NP-hard

زمان چندجمله‌ای پاسخ دقیقی برای آن‌ها به دست آورد، (مگر $P = NP$). در نتیجه، عموماً پژوهشگران این حوزه، به دنبال ارائه الگوریتم‌هایی برای حل تقریبی مسئله و یا نسخه‌های خاصی از آن بوده‌اند. یک الگوریتم تقریبی با تقریب p برای یک مسئله کمینه‌سازی، یک الگوریتم در زمان چندجمله‌ای است که یک جواب با هزینه حداکثر p برابر هزینه جواب بهینه را تضمین می‌کند. به‌طور خاص برای نسخه اصلی از مسئله k -مرکز بهترین الگوریتم ممکن که یک الگوریتم ۲-تقریب است در ۱۹۸۵ توسط گنزالس [۴] ارائه شده است.

همان‌طور که پیش‌تر بیان شد، مسئله k -مرکز به شکل‌های مختلفی تعمیم پیدا کرده است که سه نسخه ظرفیت‌دار زیر، به‌طور ویژه مورد توجه ما خواهند بود:

۱. مسئله k -مرکز با ظرفیت‌های نرم

۲. مسئله k -مرکز با ظرفیت‌های سخت

۳. نسخه وزن‌دار مسئله k -مرکز با ظرفیت‌های نرم

تاکنون نسخه‌های ظرفیت‌دار با تعاریف متعددی در پژوهش‌های مختلف مورد مطالعه قرار گرفته‌اند. خولر و سوسمان [۲] دو نسخه ظرفیت‌دار پایه‌ای از این مسئله، k -مرکز با ظرفیت‌های نرم و سخت را مورد بررسی قرار داده‌اند و دو الگوریتم با ضرایب تقریب ۵ (برای نسخه نرم) و ۶ (برای نسخه سخت) ارائه شده است.

همچنین خولر و سایرین [۵] و هونگ چان آن و همکاران [۶] نیز اخیراً بر روی نسخه‌ای کلی‌تری از مسئله که در آن ظرفیت نقاط یکسان نباشد، مطالعه انجام داده‌اند که نتیجه آن، یک الگوریتم ۹-تقریب با استفاده از شیوه برنامه‌ریزی خطی است. این نسخه از مسئله که فرم کلی‌تر نسخه‌های پیشین است، با نام « k -مرکز با ظرفیت‌های غیر یکسان^{۱۱}» شناخته می‌شود.

همان‌طور که در ابتدا بیان شد، ما به دسته دومی از پژوهش‌ها که به بحث توزیع‌شدگی پرداخته‌اند نیز توجه خواهیم داشت. مدل نگاشت-کاهش^{۱۲}، به‌عنوان یکی از چارچوب‌های^{۱۳} استاندارد محاسبات توزیع‌شده، به‌طور ویژه در بحث توزیع‌شدگی مورد توجه پژوهشگران قرار دارد و بسیاری از پژوهش‌های

^{۱۱} k -Center with non-uniform capacities

^{۱۲} MapReduce

^{۱۳} Framework

این دسته به این مفهوم پرداخته‌اند. در این مدل داده بین تعدادی ماشین توزیع می‌گردد و در هر ماشین محاسبات بر روی بخشی از داده‌ها صورت می‌پذیرد، و این عمل در طی چندین دور^{۱۴} تکرار می‌شود. یک مدل رسمی از مدل نگاشت-کاهش در پژوهش [۷] ارائه شده است.

به کارگیری این مدل برای حل مسئله k -مرکز اولین بار در پژوهش [۸] انجام شده است. همچنین در کار اخیر ایم و موسلی [۹]، مسئله k -مرکز با نقاط دورافتاده^{۱۵} در مدل نگاشت-کاهش موردبازنگری قرار گرفته است. در این پژوهش یک الگوریتم ۲-تقریب برای نسخه اصلی مسئله و یک الگوریتم ۴-تقریب برای نسخه با نقاط دورافتاده ارائه گشته است.

عموماً در مدل توزیع شده ما به دنبال آن هستیم که الگوریتم به‌طور موازی (در هر ماشین) بر روی بخشی از ورودی‌های مسئله اجرا شود. این نیاز در مفهومی بنام «مجموعه‌های هسته‌ی ترکیب‌پذیر» موردتوجه قرار گرفته است. به زبان ساده، یک مجموعه‌ی هسته‌ی ترکیب‌پذیر برای مجموعه S ، زیرمجموعه‌ای مانند $T \subseteq S$ است که حل مسئله بر روی T به ما تقریبی از جواب بهینه برای کل ورودی می‌دهد. همچنین ترکیب‌پذیری این مجموعه‌های هسته این امکان را به ما می‌دهد تا در مدلی مانند نگاشت-کاهش، با یک پردازش نهایی بر روی خروجی مرحله اول ماشین‌ها، به یک تقریب خوب از پاسخ بهینه به ازای کل داده‌های ورودی برسیم. این مفهوم در کارهای اخیر ضرابی‌زاده و سایرین [۱۰] و ایندیک و سایرین [۱۱] قابل مشاهده است.

استفاده از مفهوم مجموعه‌های هسته‌ی ترکیب‌پذیر برای حل مسائل خوشه‌بندی بسیار محدود بوده است و ما به دنبال تمرکز بر روی این حوزه هستیم. به‌طور خاص در کار اخیر باطنی و همکاران [۱۲] چند الگوریتم تقریبی با ضریب تقریب ثابت برای مسائل خوشه‌بندی و خوشه‌بندی ظرفیت‌دار در حالت توزیع شده ارائه شده است که از مفهوم مجموعه‌های هسته‌ی قابل ترکیب استفاده شده است.

۴-۱ اهداف تحقیق

در این پایان‌نامه سعی می‌شود که مسئله‌ی خوشه‌بندی نقاط، نسخه‌های مختلف آن و حل مسئله در فضای توزیع شده مورد مطالعه قرار گیرد. به‌طور ویژه مسئله k -مرکز با ظرفیت‌های نرم، k -مرکز با ظرفیت‌های

^{۱۴}Round

^{۱۵}Outlier

سخت و نسخه وزن دار آن مورد توجه ما قرار خواهد گرفت. در این پژوهش، پس از مطالعه دقیق کارهای مرتبط سعی شده است تا الگوریتم‌های موجود بازبینی گردند و در صورت امکان الگوریتم‌های جدید و بهبودیافته‌ای ارائه گردد. همچنین در عمل، در این پژوهش چهار الگوریتم جدید برای حل نسخه‌های مختلف مسئله‌ی k -مرکز ظرفیت دار ارائه گشته است.

۵-۱ ساختار پایان‌نامه

این پایان‌نامه شامل شش فصل است. فصل دوم دربرگیرنده تعاریف اولیه مرتبط با پایان‌نامه خواهد بود. در فصل سوم پژوهش‌های مرتبط با مسائل خوشه‌بندی و الگوریتم‌های توزیع شده برای حل این مسائل به تفصیل بیان می‌شود. در فصل چهارم و پنجم نتایج جدیدی که در این پایان‌نامه به دست آمده ارائه می‌گردد. در فصل چهارم، یک الگوریتم جدید برای حل نسخه وزن دار مسئله‌ی k -مرکز با ظرفیت‌های نرم ارائه می‌گردد. سپس در فصل پنجم سه الگوریتم دیگر برای حل نسخه ظرفیت دار مسئله‌ی k -مرکز ارائه می‌گردد. همچنین در پایان، بین نتایج حاصل از این پژوهش و نتایج موجود، مقایسه‌ای صورت می‌گیرد. در پایان، فصل پنجم به نتیجه‌گیری و پیشنهادهایی برای کارهای آتی خواهد پرداخت.

فصل ۲

مفاهیم اولیه

در این فصل از پایان‌نامه به تعریف و بررسی مفاهیم اولیه و پیش‌نیاز مطالبی خواهیم پرداخت که در فصل‌های آتی از آن‌ها استفاده خواهد شد.

۱-۲ الگوریتم‌های تقریبی

برای حل مسائلی که در کلاس ان‌پی-سخت قرار دارند انتظار نمی‌رود که یک الگوریتم چندجمله‌ای وجود داشته باشد. همچنین از سوی دیگر، مسئله‌ی برابری یا عدم برابری P و NP مدت‌ها است که یک مسئله‌ی حل‌نشده باقی‌مانده است. لذا برای حل این دسته از مسائل، نیازمند الگوریتم‌هایی با پیچیدگی زمانی چندجمله‌ای هستیم تا بتوانند دست‌کم جوابی نزدیک به جواب بهینه داشته باشند و درعین حال کران بالایی از حداکثر نسبت جواب تقریبی به جواب بهینه داشته باشند. لذا الگوریتم‌های تقریبی از الگوریتم‌های ابتکاری^۱ که تنها جواب نسبتاً خوبی ارائه می‌دهند تفاوت داشته و برخلاف آن‌ها می‌توانند تضمینی نسبت به کارایی جواب داشته باشند.

تعریف ۱-۲ (الگوریتم‌های تقریبی) به الگوریتم‌هایی گفته می‌شود که کران بالایی از نسبت جواب تقریبی به جواب بهینه را دارا هستند. این کران بالا ضریب تقریب^۲ خوانده می‌شود.

^۱Heuristic

^۲Approximation Factor

تعریف ۲-۲ (الگوریتم‌های تقریبی ضریب ثابت) الگوریتم‌های تقریبی که ضریب تقریب ثابت و مستقل از اندازه و شرایط مسئله‌ی ورودی دارند را الگوریتم‌های تقریبی ضریب ثابت^۳ می‌خوانیم.

یکی از مسائل کلاسیک کلاس NP ، مسئله‌ی پوشش رأسی^۴ است که الگوریتم تقریبی با ضریب ثابت ۲ برای آن وجود دارد که به صورت زیر بیان می‌گردد.

تعریف ۳-۲ (مسئله‌ی پوشش رأسی) با داشتن گراف ساده و بی‌جهت $G = (V, E)$ ، می‌خواهیم مجموعه $S \subset V$ را به گونه‌ای پیدا کنیم که $|S|$ کمینه بوده و رابطه‌ی $\forall (u, v) \in E : u \in S \vee v \in S$ برقرار باشد.

مسئله پوشش رأسی توسط یک الگوریتم حریصانه ساده حل می‌شود و به راحتی می‌توان نشان داد که این الگوریتم حریصانه، الگوریتمی تقریبی با ضریب ثابت ۲ است. برای حل مسئله به روش حریصانه، یک یال دلخواه uv را انتخاب می‌کنیم و سپس هر دو رأس u و v را به مجموعه M اضافه و کلیه یال‌های متصل به آن‌ها را از مجموعه‌ی یال‌های E حذف می‌کنیم. این کار را تا زمانی ادامه می‌دهیم که مجموعه‌ی E تهی شود. مجموعه M یک پوشش رأسی است، زیرا با تهی شدن E ، حداقل یکی از گره‌های ابتدا و انتهای یال‌ها در مجموعه‌ی M وجود خواهند داشت. همچنین مشخص است که از هر یال uv که در حلقه‌ی الگوریتم حریصانه از مجموعه E انتخاب کرده بوده‌ایم، رأس u یا v و یا هر دو، باید در مجموعه جواب بهینه موجود باشند. از آنجایی که در بدترین حالت در هر مرحله تنها یکی از رئوس مجموعه جواب بهینه را انتخاب می‌کنیم، کران بالای ضریب تقریب این الگوریتم، برابر با عدد ثابت ۲ خواهد بود.

گروه دیگر از الگوریتم‌های تقریبی، الگوریتم‌های $PTAS$ ^۵ هستند. الگوریتم‌های تقریبی از این گروه، با داشتن هر $\epsilon > 0$ به عنوان ورودی، می‌تواند جوابی با تقریب $1 + \epsilon$ (در مسئله‌های کمینه‌سازی) و یا $1 - \epsilon$ (در مسئله‌های بیشینه‌سازی) از جواب بهینه ارائه دهند. الگوریتم‌های $PTAS$ باید به ازای هر ϵ ثابت، نسبت به اندازه‌ی ورودی چندجمله‌ای باشند. در نتیجه مرتبه‌ی زمانی این الگوریتم‌ها می‌تواند حتی $O(n^{(\frac{1}{\epsilon})^c})$ باشد؛ بنابراین الگوریتم‌های $PTAS$ نسبت به ϵ ممکن است چندجمله‌ای نباشند.

در نتیجه گروه الگوریتم‌های $EPTAS$ ^۶، زیرمجموعه‌ای از گروه الگوریتم‌های $PTAS$ ، تعریف می‌شوند که زمان اجرای آن‌ها به صورت $O(n^c)$ بوده، که در این رابطه c مستقل از ϵ است. البته حتی در

^۳ Constant Factor Approximation Algorithms

^۴ Vertex Cover

^۵ Polynomial-time Approximation Scheme

^۶ Efficient Polynomial-time Approximation Scheme

الگوریتم‌های گروه $EPTAS$ نیز همچنان ϵ می‌تواند به صورت نمایی در ضریب ثابت O ظاهر شود. گروه الگوریتم‌های $FPTAS$ ^۷ امکان نمایی شدن مرتبه زمانی نسبت به اندازه مسئله و $\frac{1}{\epsilon}$ را از میان برده و تمامی الگوریتم‌های طبقه‌بندی شده در این گروه، نسبت به هر دو پارامتر اندازه مسئله و $\frac{1}{\epsilon}$ چندجمله‌ای هستند.

۲-۲ خوشه‌بندی

مسائل خوشه‌بندی یکی از قدیمی‌ترین و پایه‌ای‌ترین مسائل علوم کامپیوتر می‌باشند و در طول زمان به طور گسترده‌ای مورد مطالعه قرار گرفته‌اند. به ازای یک مجموعه نقاط ورودی، این مسائل به دنبال ارائه یک دسته‌بندی از نقاط هستند که به هر دسته در اصطلاح یک «خوشه»^۸ گفته می‌شود. همچنین معمولاً یک نقطه از هر خوشه به عنوان نماینده آن خوشه در نظر و به عنوان «مرکز خوشه»^۹ شناخته می‌شود و در نتیجه خوشه‌بندی می‌تواند تنها با مشخص کردن این مراکز صورت پذیرد.

مسائل خوشه‌بندی معمولاً به دنبال ارائه یک خوشه‌بندی همراه با کمینه‌سازی یک پارامتر خاص از آن (مانند تعداد خوشه‌ها) هستند که در اصطلاح «تابع هدف»^{۱۰} مسئله خوانده می‌شود. به همین دلیل، مسائل خوشه‌بندی به عنوان مسائل بهینه‌سازی شناخته می‌شوند. مسائلی از قبیل k -مرکز^{۱۱}، k -میانه^{۱۲} و k -میانگین^{۱۳} از جمله این مسائل خوشه‌بندی می‌باشند که در تعریف تابع هدف تفاوت‌های جزئی با یکدیگر دارند.

مسئله k -مرکز (که با نام مکان‌یابی شبکه^{۱۴} نیز شناخته می‌شود) به دنبال ارائه یک خوشه‌بندی است به طوری که بیشترین فاصله از هر نقطه تا مرکز منتسب شده به آن، که به شعاع معروف است، کمینه باشد.

^۷Fully Polynomial-time Approximation Scheme

^۸Cluster

^۹Center

^{۱۰}Objective Function

^{۱۱} k -center

^{۱۲} k -median

^{۱۳} k -means

^{۱۴}Network Location Problem

تعریف ۲-۴ (مسئله k -مرکز) مسئله k -مرکز: به ازای گراف وزن دار $G = (V, E)$ ، زیرمجموعه‌ی $S \subseteq V$ با اندازه حداکثر k پیدا کنید طوری که عبارت زیر کمینه شود.

$$\max_{u \in V} \min_{v \in S} d(u, v)$$

به‌طور مشابه، تابع هدف سایر مسائل خوشه‌بندی قابل‌تعریف خواهد بود که تعدادی از مهم‌تری آن‌ها در جدول ۲-۱ آورده شده است. همچنین در بعضی از تعاریف، توابع هدف، یا تابع هزینه^{۱۵} خوشه‌بندی، به ازای تمام مسائل خوشه‌بندی به‌صورت کلی و یکپارچه مطرح می‌شود.

تعریف ۲-۵ (تابع هزینه مسئله خوشه‌بندی p) اگر مجموعه نقاط ورودی مسئله خوشه‌بندی را V بنامیم و زیرمجموعه‌های C_1 تا C_r یک خوشه‌بندی از این نقاط باشد، به طوری که v_i به مرکز خوشه C_i اشاره کند، تابع هزینه خوشه‌بندی را می‌توان به‌صورت زیر تعریف نمود:

$$Cost(C_p) = \left(\sum_i \sum_{c \in C_i} d(v, v_i)^p \right)^{\frac{1}{p}}$$

به ازای $p = 1$ ، $p = 2$ و $p = \infty$ در تعریف ۲-۵، به توابع هدف مسائل k -میانه، k -میانگین و k -مرکز (مطابق با جدول ۲-۱) خواهیم رسید.

مسئله	تابع هدف
k -مرکز	بیشترین فاصله از هر نقطه تا مرکز منتسب شده به آن
k -میانه	میانگین فاصله نقاط یک خوشه از مرکزش
k -میانگین	مربع میانگین فاصله نقاط یک خوشه از مرکزش

جدول ۲-۱: تابع هدف تعدادی از مسائل خوشه‌بندی

از آنجا که تمام مسائل خوشه‌بندی ان‌پی-سخت^{۱۶} هستند [۱] و نمی‌توان در زمان چندجمله‌ای پاسخ دقیقی برای آن‌ها به دست آورد، (مگر $P = NP$)، عموماً پژوهشگران به دنبال ارائه الگوریتم‌هایی برای حل تقریبی مسئله بوده‌اند. به‌عنوان نمونه، بهترین الگوریتم تقریبی موجود برای مسئله k -مرکز یک الگوریتم ۲-تقریب است که در سال ۱۹۸۵ توسط گنزالس [۴] ارائه شده است.

^{۱۵} Cost Function

^{۱۶} NP-hard

مسائل خوشه‌بندی، و به‌ویژه مسئله k -مرکز، به شکل‌های مختلف قابل‌تعمیم است. در دسته‌ای از مسائل به ازای هر نقطه از ورودی به‌عنوان یک مرکز، یک هزینه بازگشایی برای آن نقطه در نظر گرفته می‌شود. در این نسخه از مسئله خوشه‌بندی که به «مسئله مکان‌یابی تسهیلات^{۱۷}» شناخته می‌شود، علاوه بر شعاع خوشه، ما به دنبال کاهش هزینه‌های بازگشایی مراکز نیز هستیم.

در دسته‌ای دیگر از پژوهش‌ها، هدف کم کردن تعداد مراکز با یک شعاع ثابت است. در این مسائل شعاع خوشه، به‌عنوان ورودی مسئله داده می‌شود و تابع هدف بر روی تعداد مراکز تعریف می‌شود. مصداق عملی این مسئله، قرار دادن تعدادی آنتن با بُرد مشخص برای پوشش تعدادی مشتری است. همچنین بخشی از مطالعات بر روی نسخه‌های «با نقاط دورافتاده^{۱۸}» از مسئله خوشه‌بندی متمرکز شده است. در این نسخه‌ها این امکان وجود دارد که تعداد محدودی از نقاط ورودی در جواب نهایی در نظر گرفته نشود و یک جواب (مستقل از وجود چنین نقاطی) برای مسئله به دست آید.

همچنین دسته دیگری از پژوهش‌ها به دنبال اعمال محدودیت بر اندازه خوشه‌ها و ایجاد توازن بین آن‌ها هستند که به مسائل خوشه‌بندی «ظرفیت‌دار^{۱۹}» مشهور هستند. در این مسائل، به هر نقطه یک ظرفیت^{۲۰} تخصیص داده می‌شود (به‌عنوان بخشی از ورودی مسئله) که نشان می‌دهد یک نقطه در صورتی که به‌عنوان مرکز انتخاب گردد تا چه تعداد نقطه را می‌تواند به خود منتسب کند. از دیدگاه ظرفیت نقاط، مسائل ظرفیت‌دار، به دو دسته «با ظرفیت‌های یکسان^{۲۱}» و «با ظرفیت‌های غیر یکسان^{۲۲}» قابل‌طبقه‌بندی است.

علاوه بر تقسیم‌بندی فوق، مسائل ظرفیت‌دار با تعاریف متعدد دیگری مانند دو مسئله « k -مرکز با ظرفیت‌های نرم^{۲۳}» و « k -مرکز با ظرفیت‌های سخت^{۲۴}» نیز در پژوهش‌های مختلف مورد مطالعه قرار گرفته‌اند. در تعریف مسئله k -مرکز با ظرفیت‌های سخت، ظرفیت تمام خوشه‌ها یکسان و برابر L در نظر گرفته می‌شود. در مسئله k -مرکز با ظرفیت‌های نرم، در کنار این محدودیت بر روی اندازه

Facility Location^{۱۷}Outlier^{۱۸}Capacitated^{۱۹}Capacity^{۲۰}Uniform^{۲۱}Non-uniform^{۲۲} k -center with soft capacities OR Capacitated multi- k -center^{۲۳} k -center with hard capacities / L-Balanced k -center^{۲۴}

خوشه‌ها، این آزادی عمل داده می‌شود تا یک نقطه بیش از یک بار بتواند به عنوان مرکز استفاده شود. در این صورت ما می‌توانیم در هر نوبت از الگوریتم، بدون در نظر گرفتن انتخاب‌های قبلی، هر یک از نقاط ورودی را به عنوان یک مرکز جدید انتخاب نماییم.

همچنین در نسخه‌ی خاصی از مسئله‌ی k -مرکز ظرفیت‌دار، به ازای هر نقطه، علاوه بر ظرفیت، یک وزن نیز قائل می‌شویم. در این صورت انتساب یک نقطه به یک مرکز، به اندازه وزن آن نقطه از ظرفیت مرکز منتسب شده استفاده خواهد کرد. این مسئله به « k -مرکز وزن‌دار با ظرفیت‌های نرم/سخت^{۲۵}» مشهور است.

۳-۲ مجموعه هسته

به ازای مجموعه نقاط P ، مجموعه‌ی $C \subset P$ را یک «مجموعه‌ی هسته^{۲۶}» می‌خوانیم به شرطی که علاوه بر $|C| < |P|$ ، محاسبه جواب بهینه بر روی مجموعه C تقریبی برای جواب بهینه بر روی مجموعه اصلی P باشد. یک مجموعه هسته «مجموعه هسته α -تقریب^{۲۷}» نامیده می‌شود اگر جواب بهینه بر روی مجموعه هسته، یک جواب α تقریب از جواب بهینه بر روی مجموعه اصلی باشد.

مفهوم مجموعه‌های هسته، به عنوان یک شیوه کارا در مدل توزیع‌شده و جریانی^{۲۸}، اولین بار توسط آگاروال و سایرین [۱۳] در سال ۲۰۰۴ استفاده شد. این مفهوم به ما کمک می‌کند تا به ازای داده‌های حجیم، مسئله بر روی بخش‌های کوچک (که بر روی یک ماشین قابل اجرا است) حل شود، و از اجتماع نتایج برای ارائه یک جواب تقریبی استفاده کنیم. این مفهوم در پژوهش‌های اخیر به عنوان یک ابزار ارزشمند برای طراحی الگوریتم‌های تقریبی (و به‌طور خاص به ازای داده‌های حجیم) مورد استقبال قرار گرفته است.

مجموعه‌های هسته به دو دسته «ترکیب‌پذیر^{۲۹}» و «تجزیه‌پذیر^{۳۰}» قابل تقسیم‌بندی می‌باشند. برای

^{۲۵} k -center with soft/hard capacities and demands

^{۲۶} Coreset

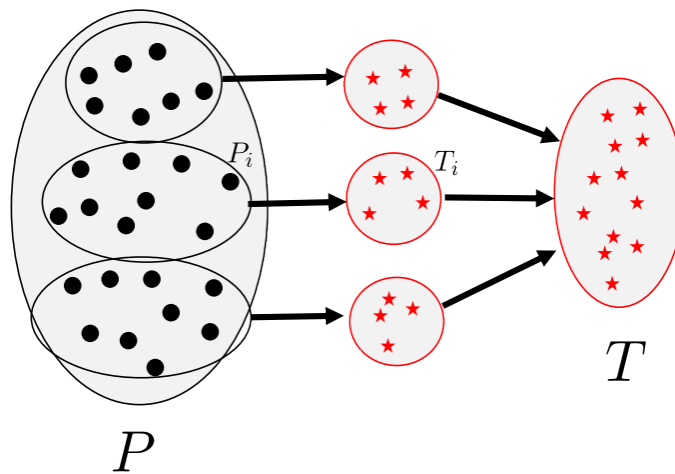
^{۲۷} α -coreset

^{۲۸} Streaming

^{۲۹} Composable

^{۳۰} Decomposable

بیان تفاوت این دو، از یک مثال استفاده می‌کنیم. در شکل ۲-۱، مجموعه ورودی P به زیرمجموعه‌های P_1 تا P_r افزار گشته است و T_i مجموعه هسته‌ی α -تقریب برای P_i است. با توجه به تعریف مجموعه‌ی هسته، این بدان معناست که جواب بهینه بر روی مجموعه T_i ، یک جواب α تقریب از جواب بهینه بر روی P_i خواهد بود. حال مجموعه T را از اجتماع T_i ها را در نظر می‌گیریم. اگر بتوان اثبات کرد که T نیز یک جواب α تقریب برای مجموعه P است، آنگاه مجموعه هسته را تجزیه‌پذیر می‌خوانیم؛ اما اگر T یک جواب با تقریب $\beta \geq \alpha$ از جواب بهینه بر روی P باشد، مجموعه هسته را ترکیب‌پذیر می‌خوانیم. طبیعی است که هر مجموعه هسته‌ی تجزیه‌پذیر، حتماً ترکیب‌پذیر نیز است.



شکل ۲-۱: مثالی از مجموعه هسته

همچنین در پژوهش اخیر باطنی و سایرین [۱۲] مفهوم جدیدی بنام «مجموعه‌ی هسته‌ی ترکیب‌پذیر با نداشت^{۳۱}» معرفی شده است. در این تعریف، علاوه بر نگهداری نقاط مجموعه‌ی هسته، یک نگاشت از نقاط مجموعه اصلی به نقاط مجموعه هسته نیز نگهداری می‌شود.

فصل ۳

کارهای پیشین

در این فصل به بررسی دقیق مجموعه پژوهش‌هایی که به نحوی با پژوهش جاری در ارتباط هستند پرداخته خواهد شد. با توجه به آنکه مسئله اصلی مورد بحث در این پژوهش ارائه یک الگوریتم تقریبی برای حل مسئله k -مرکز ظرفیت‌دار در مدل توزیع شده است، بنابراین دودسته از پژوهش‌ها مورد توجه ما خواهند بود. یک دسته پژوهش‌هایی که به مسائل خوشه‌بندی و بطور خاص مسئله k -مرکز و نسخه‌های ظرفیت‌دار آن در مدل غیرتوزیع شده پرداخته‌اند. دسته دیگر، مطالعات مربوط به مدل‌ها و شیوه‌های توزیع‌شدگی از قبیل مدل نگاشت-کاهش و پژوهش‌هایی که به نحوه‌ی حل توزیع‌شده مسائل خوشه‌بندی توجه داشته‌اند.

۳-۱ مسائل خوشه‌بندی

در این بخش به کارهای انجام‌شده به منظور حل مسئله k -مرکز و نسخه‌های ظرفیت‌دار آن خواهیم پرداخت. در این بخش ابتدا به پژوهش‌های صورت گرفته حول شکل اصلی مسئله اشاره خواهیم داشت و به تدریج بر روی پژوهش‌های مرتبط به مسئله مورد علاقه در این پژوهش، تمرکز خواهیم داشت. مسئله خوشه‌بندی به عنوان یکی از پایه‌ای‌ترین مسائل علوم کامپیوتر مورد مطالعات گسترده‌ای قرار گرفته است که در بررسی^۱ [۱۴] مروری بر این مطالعات شده است. از آنجاکه تمام مسائل خوشه‌بندی ان پی-

Survey^۱

سخت ^۲ هستند [۱] و نمی‌توان در زمان چندجمله‌ای پاسخ دقیقی برای آن‌ها به دست آورد، (مگر $P = NP$)، عموماً پژوهشگران به دنبال ارائه الگوریتم‌هایی برای حل تقریبی مسئله بوده‌اند.

یکی از اصلی‌ترین مسائل خوشه‌بندی، مسئله k -مرکز است. با توجه به آن‌پی-سخت بودن این مسئله، به‌طور خاص برای نسخه اصلی مسئله k -مرکز به ازای هیچ ϵ ای یک پاسخ تقریبی با ضریب بهتر از $\epsilon - 2$ وجود ندارد (مگر $P = NP$). تنها برای نسخه اصلی مسئله k -مرکز، الگوریتم‌ها گوناگونی ارائه شده است که در آن‌ها از شیوه‌های متفاوتی (مانند برنامه‌ریزی خطی ^۳، شیوه اولیه-دوگان ^۴ و جستجوی محلی ^۵) استفاده می‌شود.

گنزلس [۴] بهترین الگوریتم ممکن را برای نسخه اصلی از این مسئله را در سال ۱۹۸۵ ارائه داده است. این الگوریتم کار خود را با انتخاب یک نقطه دلخواه شروع می‌کند و سعی می‌کند تا $k - 1$ مرحله بعد، نقطه‌ای از میان نقاط باقی‌مانده انتخاب کند که فاصله‌اش از مجموعه نقاط انتخاب‌شده بیشینه باشد. الگوریتم بیان‌شده یک جواب 2 -تقریب ارائه می‌دهد که با توجه به توضیحات پیشین، بهترین جواب ممکن است.

همان‌طور که پیش‌تر نسخه‌های ظرفیت‌دار مسئله k -مرکز معرفی شد، نسخه‌های مختلفی از مسئله k -مرکز ظرفیت‌داری وجود دارد. با وجود مطالعات گوناگون برای شناخت نسخه‌های مختلف مسئله k -مرکز، الگوریتم‌های محدودی با ضریب تقریب ثابت برای نسخه‌های ظرفیت‌دار این مسئله وجود دارد و تا چند سال اخیر شیوه‌هایی مانند برنامه‌ریزی خطی برای حل این دسته از مسائل وجود نداشته است.

بارلان و سایرین [۱۵] در سال ۱۹۹۳ یک الگوریتم تقریبی برای حل مسئله k -مرکز ظرفیت‌دار ارائه دادند. این الگوریتم که برای نسخه ظرفیت سخت مسئله است، یک جواب 10 -تقریب ارائه می‌دهد. این ضریب تقریب توسط پژوهش خولر و سوسمان [۲] در سال ۲۰۰۰ بهبود یافت. در این پژوهش، ضریب تقریب 10 ارائه‌شده برای مسئله k -مرکز با ظرفیت‌های سخت، توسط یک الگوریتم 6 -تقریب بهبودیافته است. همچنین در این پژوهش، یک الگوریتم با ضریب تقریب 5 برای مسئله k -مرکز با ظرفیت‌های نرم ارائه شده است.

NP-hard^۲

Linear Programming^۳

Primal-Dual^۴

Local Search^۵

به صورت سطح بالا، در هر دو الگوریتم ارائه شده توسط خولر و سوسمان برای نسخه های سخت و نرم، عمل انتخاب مرکز خوشه ها را در دو گام اصلی صورت می پذیرد و این دو الگوریتم تنها در گام دوم با یکدیگر اختلاف دارند. همچنین ایده اصلی استفاده شده در گام اول هر دو الگوریتم، به کارگیری الگوریتم «شار بیشینه با هزینه کمینه»^۶ است. در الگوریتم مربوط به نسخه ظرفیت سخت مسئله، به دلیل محدودیت در نحوه انتخاب مراکز، نیاز است تا عمل انتخاب مراکز با دقت بیشتری انجام شود که این موضوع نقطه اختلاف دو الگوریتم ارائه شده در این پژوهش است.

همچنین در سال های اخیر مطالعاتی بر روی مسئله k -مرکز با ظرفیت های غیر یکسان صورت گرفته است. خولر و سایرین [۵]، برای مسئله k -مرکز با ظرفیت های غیر یکسان، اولین الگوریتم با ضریب تقریب ثابت^۷ را در سال ۲۰۱۲ ارائه کردند. این الگوریتم از شیوه رندسازی برنامه ریزی خطی^۸ استفاده می کند. همچنین خولر و سایرین در این پژوهش اثبات کرده اند که به ازای هیچ $\epsilon > 0$ ، الگوریتمی با ضریب تقریب $3 - \epsilon$ برای مسئله k -مرکز با ظرفیت های غیر یکسان وجود ندارد (مگر $P = NP$).

همچنین در این پژوهش یک الگوریتم با ضریب تقریب ۱۱ برای مسئله k -مرکز با ظرفیت های غیر یکسان برای حالتی که ظرفیت ها نرم در نظر گرفته شود ارائه شده است. هونگ چان آن و سایرین [۶] با بهبود چشمگیر پژوهش فوق، موفق شدند تا یک الگوریتم با ضریب تقریب ۹ به کمک شیوه های برنامه ریزی خطی ارائه دهند. بگفته ی مقاله، هدف بعدی این پژوهشگران، ارائه یک الگوریتم به شیوه مشابه برای مسئله مکان یابی تسهیلات است که در آن هزینه بازگشایی مراکز نیز باید در نظر گرفته شود. همچنین در پژوهش جدید سیگان و کویماکا [۱۶] در سال ۲۰۱۴، نسخه ی خاصی از مسئله k -مرکز ظرفیت دار (ترکیبی از مسئله k -مرکز با ظرفیت های غیر یکسان و نسخه همراه با نقاط دورافتاده) مورد بررسی قرار گرفته است. نتیجه این پژوهش، ارائه اولین الگوریتم تقریبی با ضریب تقریب ۲۵ است. این مسئله با فرض ظرفیت های سخت صورت گرفته است و پیش بینی می شود که با در نظر گرفتن ظرفیت ها به صورت نرم، ضریب تقریب فعلی بهبود یابد.

^۶Min-cost Max-flow

^۷حدود چندصد-تقریب که مشخصاً محاسبه نشده است

^۸LP Rounding

۲-۳ مسئله‌ی k -مرکز در مدل توزیع شده

همان‌طور که اشاره شد، به دلیل رشد سریع حجم داده‌های نیازمند پردازش، توجه پژوهشگران به استفاده از مدل‌های توزیع شده و چارچوب‌هایی همچون نگاشت-کاهش افزایش یافته است. هرچند در سال‌های اخیر، مسائل خوشه‌بندی به‌طور گسترده‌ای مورد مطالعه قرار گرفته‌اند، ولی تعداد کمی از روش‌های ابداعی در این حوزه، قابلیت استفاده برای داده‌های حجیم و مدل‌های توزیع شده را دارند.

اکثر الگوریتم‌های ترتیبی که به‌طور معمول در پژوهش‌های گذشته ارائه شده‌اند، متأسفانه با بزرگ شدن حجم داده‌ها نامطلوب ظاهر می‌شوند. یکی از دلایل این اتفاق نبود حافظه کافی است. در واقع طبیعی است که هر ماشین به اندازه محدودی حافظه دارد و در صورتی که حجم داده‌ها بیش از این اندازه باشد، ماشین عملاً قادر به اجرای الگوریتم نخواهد بود. همچنین مشکل دوم عدم امکان موازی‌سازی این الگوریتم‌ها است. در واقع الگوریتم‌های موجود، از ابتدا به‌صورت ترتیبی طراحی شده‌اند و امکان توزیع داده بین ماشین‌های مختلف را امکان‌پذیر نمی‌سازد.

مدل نگاشت-کاهش^۹ یکی از چارچوب‌های^{۱۰} استاندارد محاسبات توزیع شده است و می‌تواند به‌عنوان راه‌حل برای مشکلات فوق در نظر گرفته شود. در این مدل، داده‌ها بین تعدادی ماشین توزیع می‌شوند و در هر ماشین محاسبات بر روی بخشی از داده‌ها صورت می‌پذیرد. این عمل در طی چندین دور^{۱۱} متوالی صورت می‌پذیرد. در این مدل، امکان تبادل اطلاعات بین ماشین‌ها به شکل محدود در پایان هر دور در نظر گرفته می‌شود.

یک مدل رسمی از مدل نگاشت-کاهش در پژوهش [۷] ارائه گشته است. مهم‌ترین محدودیت در این مدل آن است که تعداد ماشین‌ها و حافظه هر ماشین باید به‌صورت زیرخطی^{۱۲} از سایز ورودی باشد که یک فرض طبیعی و قابل‌پذیرش برای داده‌های بزرگ است. همچنین با هدف طراحی هرچه موازی‌تر الگوریتم‌ها در این مدل، هیچ ماشینی در طی محاسبات نباید کل داده‌ها را مشاهده کند. زمان لازم برای هر دور چندجمله‌ای در نظر گرفته می‌شود و در واقع، هدف اصلی در این مدل کاهش تعداد دورهای لازم برای انجام الگوریتم است.

^۹ MapReduce

^{۱۰} Framework

^{۱۱} Round

^{۱۲} Sublinear

الگوریتم‌های ترتیبی فعلی برای مسئله‌ی k -مرکز در مدل نگاشت-کاهش به‌طور بهینه قابل استفاده نمی‌باشند. به‌عنوان مثال الگوریتم گنزالس را می‌توان در نظر گرفت. این الگوریتم در مرحله اول یک نقطه دلخواه انتخاب می‌کند و در مراحل بعد، دورترین نقطه از نقاط انتخاب شده تاکنون را انتخاب می‌کند. این الگوریتم در مدل نگاشت-کاهش به k دور نیاز خواهد داشت. چراکه انتخاب هر نقطه، وابسته به انتخاب‌های قبلی است و امکان موازی‌سازی در این الگوریتم مشاهده نمی‌شود.

به‌کارگیری مدل نگاشت-کاهش برای حل مسائل خوشه‌بندی اولین بار در پژوهش آینه و سایرین [۸] مورد توجه قرار گرفته است. در این پژوهش دو الگوریتم تقریبی تصادفی با ضریب تقریب ثابت برای مسئله‌ی k -مرکز و k -میان‌ارائه شده است. همچنین در مقایسه انجام شده بین روش‌های مختلف خوشه‌بندی در مدل توزیع شده، الگوریتم‌های ارائه شده در این پژوهش کارایی مناسبی از خود نشان داده‌اند.

همچنین در کار اخیر ایم و موسلی [۹] در سال ۲۰۱۵، مسئله‌ی k -مرکز و نسخه‌ی با نقاط دورافتاده^{۱۳} آن در مدل نگاشت-کاهش موردبازنگری قرار گرفته است. در این پژوهش برای نسخه اصلی مسئله، یک الگوریتم ۲-تقریب با تعداد ۴ دور و برای نسخه با نقاط دورافتاده نیز، با فرض دانستن جواب بهینه، یک الگوریتم ۴-تقریب با تعداد ۳ دور ارائه گشته است.

مفهوم مجموعه‌های هسته‌ی ترکیب‌پذیر به‌عنوان یک شیوه کارا در مدل توزیع شده و جریانی^{۱۴}، به ما کمک می‌کند تا به ازای داده‌های حجیم، مسئله بر روی بخش‌های کوچک که بر روی یک ماشین قابل اجرا است حل شود و از اجتماع نتایج، برای ارائه یک جواب تقریبی استفاده کنیم. این مفهوم به‌طور کامل در پژوهش‌های ضرابی‌زاده و سایرین [۱۰] و ایندیک و سایرین [۱۱] قابل مشاهده است. هرچند مسئله‌ی موردبررسی در این دو پژوهش (بیشینه کردن تنوع^{۱۵}) متفاوت از پژوهش جاری است، اما به‌عنوان دو نمونه ارزشمند از نحوه به‌کارگیری مفهوم مجموعه‌های هسته‌ی ترکیب‌پذیر قابل بررسی هستند.

ایده اصلی در این پژوهش‌ها توزیع داده‌ی ورودی به تعدادی ماشین است که به‌طور مجزا به حل مسئله می‌پردازند. سپس نتایج هر ماشین به‌عنوان مجموعه‌ی هسته بخشی از ورودی با یکدیگر ترکیب می‌شوند و پردازش نهایی بر روی مجموعه ترکیب شده صورت می‌پذیرد. ضرابی‌زاده و سایرین با به‌کارگیری از

^{۱۳}Outlier

^{۱۴}Streaming

^{۱۵}Diversity Maximization

ویژگی‌های مجموعه‌های هسته‌ی ترکیب‌پذیر توانستند ضرایب تقریب ارائه‌شده در کار ایندیک و سایرین را به صورت چشمگیری بهبود دهند.

با وجود مزایای گسترده‌ی این مفهوم، استفاده از آن برای حل نسخه‌های ظرفیت‌دار مسئله‌ی k -مرکز بسیار محدود بوده است. همان‌طور که بیان شد، یک نمونه از به‌کارگیری از این مفهوم در این دسته از مسائل مربوط به پژوهش نسبتاً جدید باطنی و همکاران [۱۲] است. در این پژوهش چند الگوریتم تقریبی تصادفی با ضرایب تقریب ثابت برای مسائل خوشه‌بندی و خوشه‌بندی ظرفیت‌دار در حالت توزیع‌شده (به کمک تعداد ثابتی دور در مدل نگاهت-کاهش) و در حالت جریان‌ی ارائه شده است که از مفهوم مجموعه‌های هسته‌ی قابل ترکیب استفاده شده است. در این پژوهش از نوع تعمیم‌یافته‌ای از مجموعه‌های هسته‌ی ترکیب‌پذیر استفاده شده است که نه تنها زیرمجموعه‌ای به‌عنوان معرف از نقاط اصلی را شامل می‌شود، بلکه نگاهتی از نقاط مبدأ به نقاط انتخابی را نیز به همراه دارد. در این پژوهش به‌طور خاص برای مسئله‌ی k -مرکز با ظرفیت‌های سخت، یک الگوریتم (32α) -تقریب تصادفی در تعداد دور ثابتی ارائه شده است که α به ضریب تقریب حل نسخه اصلی مسئله‌ی k -مرکز اشاره دارد. در نتیجه، الگوریتم ارائه شده در این پژوهش یک جواب ۶۴-تقریب ارائه می‌دهد. الگوریتم ارائه شده در پژوهش جاری نه تنها این ضریب تقریب را بهبود می‌دهد، بلکه الگوریتمی غیرتصادفی و با تنها ۱ دور (در مدل نگاهت-کاهش) می‌باشد.

فصل ۴

الگوریتم تقریبی جدید برای k -مرکز وزن دار

در این فصل قصد داریم تا یک الگوریتم جدید برای حالت کلی‌تر مسئله‌ی k -مرکز با ظرفیت‌های نرم که به نام k -مرکز وزن دار با ظرفیت‌های نرم شناخته می‌شود ارائه دهیم. در این نسخه از مسئله برای هر نقطه یک وزن و یا یک میزان درخواست^۱ در نظر گرفته می‌شود. میزان درخواست یک نقطه می‌تواند توسط یک یا چند مرکز تأمین گردد. زمانی که یک نقطه به یک مرکز انتساب پیدا می‌کند، به اندازه وزن/درخواست آن نقطه از ظرفیت مرکز استفاده می‌شود. ارائه این الگوریتم به ما کمک می‌کند تا در فصل بعد چند الگوریتم تقریبی با ضریب تقریب ثابت برای مسئله k -مرکز با ظرفیت‌های نرم و سخت طراحی نماییم. شایان ذکر است که نتایج حاصل در این فصل و فصل بعد با همکاری جناب مهندس کیان میرجلالی بدست آمده است.

الگوریتم ارائه شده در این فصل بر مبنای الگوریتم خولر و سوسمان [۲] شکل گرفته است. الگوریتم مذکور، یک الگوریتم ۵-تقریب برای نسخه غیر وزن دار مسئله k -مرکز با ظرفیت‌های نرم می‌باشد، که ما در این بخش آن را به نسخه‌ی وزن دار ارتقا می‌دهیم. به منظور اجتناب از بیان بخش‌های تغییر نکرده‌ی الگوریتم، بعضاً به قسمت‌هایی از آن ارجاع داده خواهد شد. الگوریتم نهایی یک جواب ۵-تقریب برای مسئله k -مرکز وزن دار با ظرفیت‌های نرم خواهد بود.

^۱Request

۱-۴ تعاریف و مفاهیم اولیه

برای تعریف دقیق مسئله لازم است ابتدا چند تعریف پایه‌ای ارائه دهیم. همان‌طور که پیش‌تر بیان شد، در نسخه وزن دار مسئله، هر رأس u یک میزان درخواست r_u دارد که به‌عنوان ورودی مسئله داده می‌شود. مجموع تمام درخواست رؤس را با R_{total} نمایش می‌دهیم و به‌صورت زیر تعریف می‌شود:

$$R_{total} = \sum_{v \in V} r_v$$

در نسخه وزن دار از مسئله، میزان درخواست هر رأس می‌تواند توسط چندین مرکز تأمین شود. مجموعه مراکز که رأس u از آن‌ها سرویس دریافت می‌کند و به آن‌ها منتسب شده است را با نماد $\phi(u)$ نمایش می‌دهیم.

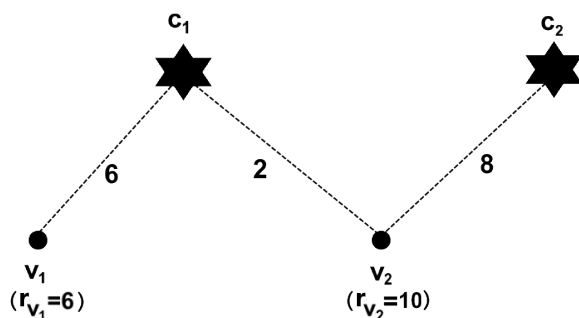
مجموعه رؤسی که از یک مرکز سرویس‌گیری می‌کنند و میزان این سرویس‌گیری، دامنه \mathcal{D} یک مرکز را تعریف می‌کند. دامنه یک مرکز مانند c در واقع از دو جز تشکیل می‌شود: (۱) مجموعه رؤسی که بخشی از درخواست خود را توسط این مرکز تأمین می‌کند. این مجموعه را با نماد D_c نمایش می‌دهیم. (۲) میزان درخواستی از هر عضو مجموعه فوق که توسط مرکز تأمین می‌گردد، که با نماد F_c و به‌صورت تابع زیر تعریف می‌شود:

$$F_c : D_c \rightarrow \mathbb{Z}$$

همچنین به ازای هر مرکز m ، «اندازه دامنه m » را با نماد R_m و برابر با میزان درخواست‌های پوشش داده‌شده توسط m تعریف می‌کنیم و داریم:

$$R_m = \sum_{v \in \mathcal{D}_m} F_m(v)$$

برای واضح شدن نمادهای معرفی شده، به مثال شکل ۴-۱ که در آن دو رأس v_1 و v_2 از دو مرکز c_1 و c_2 سرویس می‌گیرند توجه کنید.



شکل ۴-۱: یک مثال برای معرفی نمادهای معرفی شده در این بخش

$$R_{total} = 16$$

$$R_{c_1} = 8 \quad R_{c_2} = 8$$

$$D_{c_1} = \{v_1, v_2\}$$

$$F_{c_1} = \{v_1 \rightarrow 6, v_2 \rightarrow 2\}$$

$$D_{c_2} = \{v_2\}$$

$$F_{c_2} = \{v_2 \rightarrow 8\}$$

$$\phi(v_1) = \{c_1\}$$

$$\phi(v_2) = \{c_1, c_2\}$$

همچنین از پیش می دانستیم که فاصله یک نقطه p از یک مجموعه نقاط S به صورت فاصله نقطه p از نزدیک ترین نقطه مجموعه S تعریف و با نماد $d(p, S)$ نمایش داده می شود. با توجه به تعاریف فوق، به بیان تعریف رسمی مسئله می پردازیم.

مسئله k -مرکز وزن دار با ظرفیت های نرم: به ازای گراف ورودی $G = (V, E)$ که طول یال e_i با $w(e_i)$ نمایش داده شده است و میزان درخواست (یا وزن) هر رأس مانند u به شکل r_u مشخص شده است، یک زیرمجموعه $S \subseteq V$ با اندازه حداکثر k و نحوه انتساب سایر رئوس به این زیرمجموعه پیدا کنید که حداکثر فاصله هر رأس از دورترین مرکز متناسب به آن کمینه باشد. به عبارت دیگر:

$$\min_{S \subseteq V} \max_{\substack{u \in V \\ p \in \phi(u)}} d(u, p)$$

به طوری که:

$$\sum_{m \in \phi(u)} F_m(u) \leq r_u \quad \forall u \in V$$

$$\sum_{u \in \mathcal{D}_v} F_v(u) \leq L \quad \forall v \in S$$

۲-۴ الگوریتم پیشنهادی

در این بخش به بیان الگوریتم پیشنهادی خواهیم پرداخت. درستی الگوریتم در بخش بعد مورد بررسی خواهد گرفت. همان طور که پیش تر بیان شد، برای اجتناب از تکرار تا جای ممکن در بخش هایی به الگوریتم اصلی ارجاع خواهیم نمود. همچنین در ادامه فرض می شود که گراف ورودی یک گراف کامل است؛ چراکه به ازای هر یالی که وجود نداشته باشد، می توان کوتاه ترین مسیر^۳ بین آن دو رأس را در نظر گرفت و الگوریتم همچنان به درستی عمل خواهد کرد.

در بالاترین سطح از الگوریتم، قطعه کد $\text{CAPACITED-CENTERS}(G = (V, E), k, L)$ ، ابتدا یال ها را بر اساس وزن به صورت غیر نزولی مرتب می کند. یال های مرتب شده را e_1, e_2, \dots, e_n در نظر بگیرید. به ازای هر i ، گراف G_i را یک زیرگراف بی وزن از G شامل یال های با حداکثر طول $w(e_i)$ در نظر بگیرید. از آنجاکه در تمام طول الگوریتم گراف G_i بی وزن در نظر گرفته می شود، میدانیم که فاصله هر دو رأس در آن برابر با «تعداد» یال های موجود در کوتاه ترین مسیر بین آن دو رأس خواهد بود. به ازای هر گراف G_i (با شروع i از ۱ تا n)، و تا زمانی که به جواب برسیم، الگوریتم $\text{ASSIGNCENTERS}(G_i)$ را بر روی G_i اجرا می کنیم. در پایان هر مرحله تعدادی مرکز به عنوان جواب پیدا خواهد شد که دو حالت خواهد بود:

- اگر تعداد این مرکزها از k بیشتر باشد، ما اثبات خواهیم کرد که هیچ راه حلی با حداکثر هزینه $w(e_i)$ وجود ندارد. زیرا طبق قضیه ۴-۵ و ۴-۴ الگوریتم ما در بدترین حالت تنها به اندازه حداقل تعداد مراکز لازم در جواب بهینه، مرکز ایجاد می کند. پس اگر الگوریتم ما بیش از k مرکز را برای پوشش درخواست ها استفاده کند، این بدان معناست که هر جواب بهینه ای نیز به بیش از این تعداد مرکز نیاز خواهد داشت و در نتیجه جوابی وجود ندارد.
- اگر تعداد مراکز به دست آمده حداکثر k باشد، نشان خواهیم داد که در G_i ، فاصله هر رأس تا دورترین مرکزی که به آن منتسب شده حداکثر ۵ خواهد بود که برابر با یک جواب ۵-تقریب برای مسئله خواهد بود. این حالت در قضیه ۴-۵ اثبات خواهد شد.

الگوریتم ۱ $Capacitated - Centers(G(V, E), k, L)$

۱: تمام یال‌ها را به صورت غیر نزولی مرتب کن (e_1, e_2, \dots, e_n)

۲: به ازای $i = 0$ تا n انجام بده:

۳: ساخت گراف $G_i = (V, E_i)$ به طوری که $E_i = \{e_1, \dots, e_i\}$

۴: پایان الگوریتم اگر $AssignCenters(G_i) = true$

الگوریتم ۲ $AssignCenters(G_i)$

۱: قرار بده $SUCCESSFUL = true$

۲: فرض کنیم $R_i^c = \sum_{v \in c} r_v$ مجموع میزان درخواست‌های مؤلفه‌ی همبندی G_i^c باشد

۳: اگر $\sum_{c \in G_i} \left\lceil \frac{R_i^c}{L} \right\rceil > k$

۴: $SUCCESSFUL = false$

۵: در غیر این صورت:

۶: به ازای هر مؤلفه‌ی همبندی G_i^c در G_i :

۷: فراخوانی $ASSIGNMONARCHS(G_i^c)$

۸: فراخوانی $ASSIGNDOMAINS(G_i^c)$

۹: فراخوانی $REASSIGN(G_i^c)$

۱۰: اگر مجموع تعداد مراکز ایجاد شده بزرگ‌تر از k :

۱۱: $SUCCESSFUL = false$

۱۲: $SUCCESSFUL$ را برگردان

الگوریتم $AssignCenters(G_i)$ ^۴ وظیفه مشخص کردن مراکز در گراف ورودی G_i را بر عهده دارد. در این الگوریتم هر مؤلفه‌ی همبندی به صورت مجزا مورد پردازش قرار می‌گیرد. این کار از آن جهت درست است که شعاع هر جواب بهینه (در صورت وجود چنین جوابی) حداکثر $w(e_i)$ خواهد بود و در نتیجه هیچ رأسی از یک مؤلفه‌ی همبندی به یک رأس از مؤلفه‌ی دیگر منتسب نخواهد شد و در نتیجه می‌توان مؤلفه‌ی همبندی را مجزا پردازش نمود.

^۴ با اندکی تغییر نسبت به الگوریتم همتای خود در پژوهش [۲]

اگر به ازای هر مؤلفه همبندی c از گراف G_i ، داشته باشیم $R_i^c = \sum_{v \in c} r_v$ ، یک حد پایین برای تعداد مراکز لازم در این مؤلفه همبندی $\lceil \frac{R_i^c}{L} \rceil$ خواهد بود. در نتیجه اگر $\lceil \frac{R_i^c}{L} \rceil$ از حداکثر تعداد مرکز مجاز (k) بیشتر باشد، میدانیم که هیچ جوابی برای گراف G با حداکثر اندازه $w(e_i)$ وجود نخواهد داشت. در غیر این صورت، به ازای هر مؤلفه همبندی مرکزها به صورتی که در ادامه بیان می شود، به دست می آید.

در ادامه‌ی این الگوریتم و توسط فراخوانی الگوریتم $\text{ASSIGNMONARCHS}(G_i)$ (رجوع به پژوهش [۲])، به ازای هر مولفه همبندی مثل G_i^c ، یک مجموعه مستقل^۵ به نام M به دست می آید. اعضای این مجموعه مستقل که در اصطلاح پژوهش اصلی پادشاه^۶ نامیده می شوند، یک درخت (رابطه پدر-پسر) را تشکیل می دهند. هر یک از سایر رئوس مؤلفه همبندی نیز به یک پادشاه منتسب می گردند و در اصطلاح امپراتوری پادشاه را شکل می دهند. به ازای هر پادشاه مانند m ، پدر پادشاه در درخت پادشاهها با نماد $p(m)$ نمایش داده می شود. همچنین مجموعه رئوسی که طبق الگوریتم به پادشاه m منتسب شده اند و امپراتوری آن را تشکیل داده اند، با نماد $E(m)$ نمایش داده می شوند. همچنین برای سهولت اشاره به پادشاه صاحب یک امپراتوری، نماد M را به صورت زیر تعریف می کنیم.

$$\forall v \in E(m) : \mathcal{M}(v) = m$$

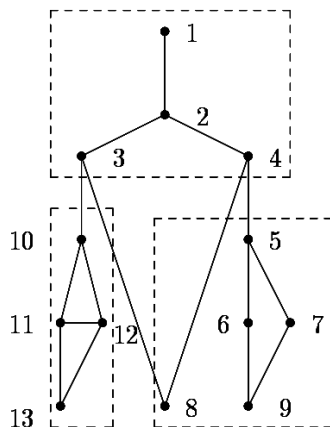
با توجه به الگوریتم مذکور، چهار ویژگی مهم زیر در رابطه با مجموعه پادشاهها و امپراتوریها وجود دارد:

۱. فاصله هر دو پادشاه از یکدیگر حداقل سه است.
۲. فاصله هر پادشاه m (به جز پادشاه ریشه درخت) تا «پادشاه پدرش» دقیقاً برابر سه است.
۳. فاصله هر پادشاه از رئوس موجود در امپراتوری خودش حداکثر دو است. رئوس در فاصله یک از پادشاه را «امپراتوری سطح اول m » یا $E^1(m)$ و رئوس در فاصله دو را «امپراتوری سطح دوم m » یا $E^2(m)$ می خوانیم و داریم: $E(m) = E^1(m) \cup E^2(m)$.
۴. هر پادشاه m (به جز پادشاه ریشه درخت) حداقل یک یال ارتباطی با یک رأس قرارگرفته در $E^2(p(m))$ دارد. همچنین با توجه به ویژگی ۱، هر چنین رأس سطح دومی، تنها یک پادشاه در مجاورت خود دارد.

Independent Set^۵

Monarch^۶

شکل ۲-۴ یک نمونه از اجرای الگوریتم $ASSIGNMONARCHS(G_i)$ بر روی یک گراف با یک مولفه‌ی همبندی و ۱۳ رأس را نمایش می‌دهد. در اینجا رئوس ۱ و ۵ و ۱۰ به‌عنوان پادشاه‌ها انتخاب شده‌اند و ناحیه نقطه‌چین، امپراتوری هر یک از این پادشاه‌ها را مشخص کرده است و داریم:



شکل ۲-۴: مثالی از مجموعه پادشاه‌ها و امپراتوری

$$M = \{1, 5, 10\}$$

$$p(5) = p(10) = 1 \quad \& \quad p(1) = \emptyset$$

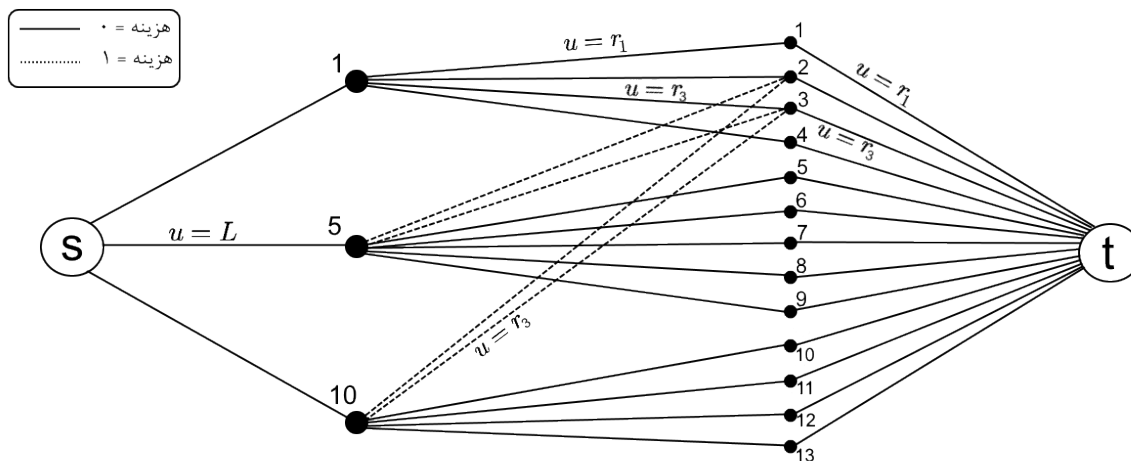
$$\mathcal{M}(2) = \mathcal{M}(3) = \mathcal{M}(4) = 1$$

پس از به دست آمدن مجموعه پادشاه‌های M ، الگوریتم $ASSIGNDOMAINS(G_i^c)$ ^۷ اجرا می‌گردد. این الگوریتم تلاش می‌کند تا یک دامنه به اندازه حداکثر L به هر پادشاه انتساب بدهد. ایده اصلی استفاده‌شده در این الگوریتم، استفاده از الگوریتم «بیشترین شار با هزینه کمینه» بر روی یک گراف دوبخشی $G'(M, V, E')$ است. به ازای گراف ورودی شکل ۲-۴ این گراف دوبخشی به صورت شکل ۳-۴ خواهد بود. هدف نهایی این الگوریتم، انتساب حداکثری درخواست رئوس به پادشاه‌ها است، به طوری که سه محدودیت زیر رعایت گردد:

۱. هر پادشاه در دامنه خود بیش از L درخواست را پوشش ندهد.
۲. یک رأس تنها به پادشاه‌هایی انتساب پیدا کند که حداکثر در فاصله دو از یکدیگر باشند.

^۷ با اندکی تغییر نسبت به الگوریتم همتای خود در پژوهش [۲]

۳. یک پادشاه تنها زمانی به درخواست یک رأس از خارج امپراتوری خود سرویس دهد که هیچ درخواست بی پاسخی در رئوس امپراتوری خود نداشته باشد.



شکل ۴-۳: مثالی از گراف دوبخشی G'

برای تطبیق این الگوریتم با نسخه وزن دار، لازم است تا تغییراتی در بدنه الگوریتم اصلی ایجاد گردد و تعدادی از تعاریف موجود گسترش یابند. از آنجاکه در نسخه وزن دار، هر رأس یک میزان درخواست مشخص (r_v) دارد، نیاز است تا ظرفیت (u) یال‌های گراف دوبخشی G' به صورت زیر اصلاح شود:

$$\forall m \in M, v \in V : \quad u(s, m) = L, \quad u(m, v) = r_v, \quad u(v, t) = r_v$$

واضح است که خروجی این الگوریتم، یک دامنه به ازای هر پادشاه خواهد بود. همان‌طور که پیش‌تر نیز بیان شد، تعریف دامنه یک مرکز در نسخه وزن دار اندکی متفاوت است و از دو جز (یک مجموعه رئوس و یک تابع) تشکیل می‌شود.

پس از اجرای الگوریتم $\text{ASSIGNDOMAINS}(G_i^c)$ ، ممکن است بخشی از درخواست‌های یک رأس توسط هیچ پادشاهی پوشش داده نشده باشد. به همین دلیل، الگوریتم $\text{REASSIGN}(G_i^c)$ ^۸ سعی می‌کند تا با پیمایش درخت پادشاه‌ها، از سمت برگ‌ها به ریشه، تعدادی مرکز جدید برای پوشش این بخش از درخواست‌ها بر روی پادشاه‌های فعلی ایجاد نماید (این کار به دلیل نرم بودن ظرفیت‌ها امکان‌پذیر است). این الگوریتم درخواست‌های مربوط به یک رأس u را به $M(u)$ یا $p(M(u))$ منسوب می‌کند.

^۸ با اندکی تغییر نسبت به الگوریتم همتای خود در پژوهش [۲]

 الگوریتم ۳ $AssignDomains(G_i^c)$

- ۱: مجموعه M را مجموعه پادشاه‌های G_i^c در نظر بگیر.
 - ۲: قرار بده $E' = \{(m, v) | m \in M, v \in V, d(m, v) \leq 2\}$.
 - ۳: گراف دوبخشی $G' = M, V, E'$ را بساز.
 - ۴: دو رأس s و t را اضافه کن.
 - ۵: یال‌های $\{(s, m) | m \in M\}$ و $\{(v, t) | v \in V\}$ را اضافه کن.
 - ۶: به ازای هر $m \in M, v \in V$ ظرفیت یال‌ها را قرار بده: $u(s, m) = L, u(m, v) = r_v, u(v, t) = r_v$.
 - ۷: هزینه هر یال $c(m, v)$ که در امپراتوری m نیست برابر ۱ و هزینه سایر یال‌ها برابر ۰ باشد.
 - ۸: الگوریتم «بیشینه شار با هزینه کمینه» را بر روی G' اجرا کن.
 - ۹: به ازای هر پادشاه $m \in M$:
 - ۱۰: قرار بده $\mathcal{D}_m = \{v | \text{رأس } v \text{ دست کم یک واحد شار از } m \text{ دریافت کرده است}\}$
 - ۱۱: $\forall v \in \mathcal{D}_m$: قرار بده $F_m(v) = \{m | \text{رأس } v \text{ از } m \text{ دریافت کرده است}\}$
 - ۱۲: به ازای هر رأس $v \in V$:
 - ۱۳: قرار بده $\phi(v) = \{m | v \in \mathcal{D}_m\}$
-

الگوریتم ۴ $ReAssign(G_i^c)$

۱: مجموعه M را مجموعه پادشاه‌های G_i^c در نظر بگیر.

۲: به ازای هر پادشاه $m \in M$:

۳: قرار بده: $Unsatisfied(m) = \left\{ v \in \{m\} \cup E(m) \mid r_v > \sum_{m \in \phi(v)} F_m(v) \right\}$

۴: قرار بده: $U(m) = \sum_{v \in Unsatisfied(m)} (r_v - \sum_{m \in \phi(v)} F_m(v))$

۵: قرار بده: $Passed(m) = \emptyset$

۶: قرار بده: $P(m) = \bullet$

۷: T را درخت پادشاه‌های G_i^c قرار بده.

۸: تا وقتی درخت T خالی نشده است:

۹: برگ m از T را در نظر بگیر.

۱۰: مرکز واقع شده بر روی پادشاه m را c بنام. (در الگوریتم $AssignDomains$)

۱۱: مقدار k' و ϵ را از رابطه زیر به دست بیاور:

$$U(m) + P(m) = k'.L + \epsilon$$

۱۲: تعداد k' مرکز جدید (بنام‌های $(c_1, \dots, c_{k'})$ بر روی رأس m ایجاد کن و $k'.L$ واحد از

درخواست‌ها را به آن‌ها منتسب کن. به ازای هر رأس v که بخشی از درخواستش توسط این

مراکز جدید پوشش داده می‌شود، رأس m را به مجموعه $\phi(v)$ اضافه کن.

۱۳: مقدار ϵ واحد درخواست باقیمانده را c منتسب کن. (ϕ مجدداً به‌روزرسانی شود)

۱۴: میزان ϵ واحد از درخواست‌هایی که مرکز c توسط الگوریتم $AssignDomains$ پوشش داده

بوده را به $Passed(p(m))$ اضافه کن.

۱۵: برگ m از T را پاک کن.

۳-۴ اثبات درستی الگوریتم

به سادگی می‌توان دید که پیچیدگی الگوریتم ما از مرتبه چندجمله‌ای است. حال در این بخش به اثبات درستی الگوریتم و محاسبه ضریب تقریب آن می‌پردازیم. بدین منظور ابتدا با توجه به الگوریتم $\text{ASSIGNDOMAINS}(G_i^c)$ تعدادی تعریف اولیه ارائه می‌دهیم. با توجه به این تعاریف به یک افراز بندی از مجموعه پادشاه‌ها خواهیم رسید. در ادامه تعدادی ویژگی ارائه خواهد شد و نهایتاً در قضیه ۴-۵ به اثبات نهایی درستی و ضریب تقریب الگوریتم پرداخته خواهد شد.

تعریف ۴-۱ (شار ورودی و خروجی به یک رأس) با توجه به اجرای الگوریتم شار بیشینه با هزینه کمینه بر روی گراف دوبخشی G' داریم:

$$\forall v \in M \cup V : f^-(v) = \left\{ \text{مجموع شار وارد شده به رأس } v \text{ پس از اجرای الگوریتم} \right\}$$

$$\forall v \in M \cup V : f^+(v) = \left\{ \text{مجموع شار خارج شده به رأس } v \text{ پس از اجرای الگوریتم} \right\}$$

تعریف ۴-۲ (رأس پوشش داده نشده) رأس v را پوشش داده نشده گوئیم اگر $f^+(v) < r_v$.

تعریف ۴-۳ با توجه شار ورودی (f^-) و خروجی (f^+) رئوس گراف دوبخشی G' ، پادشاه‌ها را به سه دسته زیر می‌توان افراز نمود.

۱. پادشاه m را سبک (**Light**) گوئیم اگر اندازه دامنه (R_m) آن کمتر از L باشد، یا به عبارت دیگر:

$$f^-(m) < L$$

۲. پادشاه m را سنگین (**Heavy**) گوئیم اگر یک رأس پوشش داده نشده در امپراتوری m وجود داشته باشد، یا به عبارت دیگر:

$$f^-(m) = L \quad \text{و} \quad f^+(E(m)) < \sum_{v \in E(m)} r_v$$

۳. پادشاه m را تکمیل (**Full**) گوئیم اگر نه سبک و نه سنگین باشد، یا به عبارت دیگر:

$$f^-(m) = L \quad \text{و} \quad f^+(E(m)) = \sum_{v \in E(m)} r_v$$

همچنین مجموعه تمام پادشاه‌های سبک، سنگین و تکمیل را به ترتیب با نمادهای LM و HM و FM معرفی می‌کنیم.

تعریف ۴-۴ تعداد کل پادشاه‌های سبک را K_{LM} می‌نامیم. همچنین مجموع درخواست‌های موجود در دامنه پادشاه‌های سبک را با نماد R_{LM} معرفی و به صورت زیر تعریف می‌کنیم:

$$R_{LM} = \sum_{l \in LM} R_l$$

لم ۴-۱ اگر m یک پادشاه سنگین باشد، تمام درخواست‌های دامنه m مربوط به امپراتوری m است، یا به عبارت دیگر:

$$\forall u \in \mathcal{D}_m \implies u \in E(m)$$

اثبات. فرض کنید رأس u در دامنه m قرار داشته باشد ($F_m(v) > 1$)، درحالی که در امپراتوری m قرار نداشته باشد. حال فرض کنید یک رأس دیگر مانند x در امپراتوری m وجود داشته باشد که پوشش داده نشده باشد.

در این صورت ما می‌توانیم نحوه توزیع شار را به این صورت تغییر دهیم که یک واحد شار عبوری از m که به u وارد شده را به x بدهیم. میزان کل شار جدید برابر با کل شار قبلی است ولی با هزینه کمتر خواهد بود که این در تناقض با الگوریتم «شار بیشینه با هزینه کمینه» است. \square

مجموعه اولیه \mathcal{E} را برابر LM (مجموعه پادشاه‌های سبک) در نظر بگیرید. در مرحله زام، پادشاهی را که یک رأس مانند v در دامنه خود دارد، و با توجه به گراف دوبخشی G' امکان پوشش داده شدن v توسط یکی از پادشاه‌های موجود در مجموعه \mathcal{E}_{j-1} نیز وجود داشته است، به مجموعه اضافه کنید. به عبارت دیگر:

$$\mathcal{E}_j = \mathcal{E}_j \cup \{m \in M \mid \exists v \in V, \exists m' \in \mathcal{E}_{j-1}, m \in \phi(v) \text{ and } d(v, m') \leq 2 \text{ in } G'\}$$

\mathcal{E} را بزرگ‌ترین مجموعه \mathcal{E}_j به دست آمده در فرآیند فوق بنامیم. همچنین \mathcal{F} را سایر پادشاه‌های قرار نگرفته در این مجموعه می‌نامیم ($\mathcal{F} = M - \mathcal{E}$).

لم ۲-۴ مجموعه \mathcal{E} شامل هیچ پادشاه سنگینی نیست.

اثبات. اثبات مشابه با اثبات ارائه شده در پژوهش [۲] است. فرض کنیم که پادشاه سنگین θ در مرحله v به واسطه رأس مشترک v به \mathcal{E}_{j-1} اضافه شده باشد. ما می توانیم یک واحد درخواست رأس v را از θ آزاد کرده و به یک پادشاه θ' منتقل کنیم. با تکرار یک دنباله از این عمل به \mathcal{E} خواهیم رسید که میدانیم ظرفیت خالی برای یک واحد منتقل شده در این سلسله جابجایی ها را حتماً دارد. با این کار ما پادشاه سنگین θ را برای سرویس دادن به یک واحد درخواست جدید به رئوس امپراتوری θ که پوشش کامل داده نشده است آماده می کنیم و شار کل را افزایش می دهیم که در تناقض با الگوریتم «شار بیشینه با هزینه کمینه» است. \square

لم ۳-۴ یک مرکز مانند θ در جواب بهینه را در نظر بگیرید که بخشی از درخواست های یک پادشاه مانند $e \in \mathcal{E}$ را پوشش می دهد. هیچ رأسی که در دامنه پادشاه های موجود در \mathcal{E} قرار نداشته باشد، نمی تواند در جواب بهینه نیز توسط θ مورد پوشش قرار بگیرد و در دامنه θ باشد.

اثبات. فرض کنیم که مرکز بهینه θ دو رأس $e \in \mathcal{E}$ و u را در دامنه خود پوشش می دهد. به عنوان فرض خلف، فرض می کنیم که u در دامنه هیچ پادشاهی از \mathcal{E} نباشد، یا به عبارت دیگر u یکی از شرایط زیر را داشته باشد:

۱. رأس u کلاً در دامنه هیچ پادشاهی نباشد:

از آنجاکه میدانیم فاصله θ تا هر دو رأس مذکور حداکثر $w(e_i)$ است، پس متوجه می شویم که در گراف بی وزن G_i داریم: $d(e, u) \leq 2$. مشابه به اثبات لم ۲-۴، می توان با دنباله ای از جابجایی ها، یک واحد درخواست از e آزاد نمود و به یک پادشاه در \mathcal{E} منتقل نمود. بنا بر لم ۲-۴، پادشاه نهایی ظرفیت خالی دارد و این سلسله جابجایی ها ممکن است. در نتیجه می توان شار را به گونه ای تغییر داد که u توسط e پوشش داده شود که در تناقض با فرض اولیه است.

۲. رأس u در دامنه پادشاه $f \in \mathcal{F}$ باشد:

به طور مشابه میدانیم که $d(e, u) \leq 2$. با توجه به تعریف \mathcal{E} میدانیم که پادشاه f می تواند عضو \mathcal{E} شود که در تناقض با فرض اولیه است.

\square

قضیه ۴-۴ هر جواب برای مسئله k -مرکز وزن دار با ظرفیت‌های نرم به ازای گراف G_i به حداقل $K_{LM} + \left\lceil \frac{R-R_{LM}}{L} \right\rceil$ مرکز نیاز خواهد داشت.

اثبات. فرض کنیم که مجموعه درخواست‌های دامنه پادشاه‌های \mathcal{E} را $R_{\mathcal{E}} = \sum_{m \in \mathcal{E}} R_m$ بنامیم. در جواب بهینه، هر پادشاه از \mathcal{E} باید توسط یک (یا چند) مرکز متمایز پوشش داده شوند^۹ و بنا بر لم ۳-۴ می‌دانیم که این مراکز جواب بهینه هیچ پادشاه دیگری از \mathcal{F} را نیز پوشش نمی‌دهند. در نتیجه، هر جوابی حداقل به $|\mathcal{E}| + \left\lceil \frac{R-R_{\mathcal{E}}}{L} \right\rceil$ مرکز نیاز دارد. همچنین داریم:

$$|\mathcal{E}| = K_{LM} + \frac{|\mathcal{E}| - K_{LM}}{L} \cdot L \quad (۴-۱)$$

$$R_{\mathcal{E}} = R_{LM} + (|\mathcal{E}| - K_{LM}) \cdot L \quad (۴-۲)$$

پس داریم:

$$|\mathcal{E}| + \left\lceil \frac{R - R_{\mathcal{E}}}{L} \right\rceil = K_{LM} + \left\lceil \frac{R + (|\mathcal{E}| - K_{LM}) \cdot L - R_{\mathcal{E}}}{L} \right\rceil \quad (\text{طبق ۱-۴})$$

$$= K_{LM} + \left\lceil \frac{R - R_{LM}}{L} \right\rceil \quad (\text{طبق ۲-۴})$$

□ در نتیجه هر جواب به حداقل $K_{LM} + \left\lceil \frac{R-R_{LM}}{L} \right\rceil$ مرکز نیاز خواهد داشت.

قضیه ۵-۴ در خروجی الگوریتم پیشنهادی ما حداکثر فاصله هر رأس از دورترین مرکزی که از آن سرویس می‌گیرد ۵ است. همچنین حداکثر $K_{LM} + \left\lceil \frac{R-R_{LM}}{L} \right\rceil$ مرکز استفاده می‌شود.

اثبات. بخش اول:

با توجه به الگوریتم $\text{ASSIGNDOMAINS}(G_i^c)$ میدانیم که بخشی از درخواست رؤس توسط پادشاهی که از آن شار دریافت می‌کند پوشش داده شده است. همچنین به ازای آن دسته از رؤوسی که به صورت ناقص^۹ در واقع مرکزی در جواب بهینه وجود ندارد که هم‌زمان چند پادشاه موجود در \mathcal{E} را پوشش بدهد. چراکه در این صورت در واقع فاصله آن دو (یا چند) مرکز از هم دو است و طبق الگوریتم $\text{ASSIGNMONARCHS}(G_i)$ عملاً یکی از آن‌ها نمی‌تواند اصلاً پادشاه باشد.

پوشش داده شده اند (طبیعتاً متعلق به پادشاه‌های سنگین)، در الگوریتم $\text{ReAssign}(G_i^c)$ این بخش از درخواست‌ها توسط پادشاه امپراتوری خودشان، و یا حداکثر پادشاه یک سطح بالاتر در درخت، پوشش داده می‌شود. در حالت اول و طبق تعریف گراف دوبخشی G' ، حداکثر فاصله دو است. در حالت دوم نیز این فاصله حداکثر پنج خواهد بود (دو واحد تا پادشاه خودشان و سه واحد تا پادشاه پدر). آر آنجا که گراف G_i کوچک‌ترین گرافی است که دارای یک جواب بهینه است، پس $OPT \geq w(e_i)$ است. همچنین همان‌طور که بیان شد، شعاع جواب الگوریتم پیشنهادی حداکثر $5 \cdot w(e_i)$ هست. در نتیجه جواب الگوریتم، یک جواب 5 -تقریب از جواب بهینه خواهد بود:

$$OPT \geq w(e_i) \quad \& \quad r \leq 5 \cdot w(e_i)$$

$$\implies r \leq 5 \cdot OPT$$

بخش دوم:

پس از اجرای الگوریتم $\text{AssignDomains}(G_i^c)$ ، هر پادشاه (به‌جز پادشاه‌های سبک) L واحد از درخواست‌ها را پوشش داده‌اند. همچنین بدیهی است که در طی الگوریتم $\text{AssignDomains}(G_i^c)$ اندازه دامنه پادشاه‌های سبک کاهش پیدا نخواهد کرد^{۱۰}. همچنین می‌دانیم که تمام مراکز جدید دقیقاً L درخواست را تحت پوشش قرار داده‌اند. در نتیجه مجموعاً $K_{LM} + \left\lceil \frac{R-R_{LM}}{L} \right\rceil$ مرکز ایجاد می‌گردد. \square

نتیجه‌ی ۴-۶ با قبول اثبات قضیه ۴-۵، الگوریتم پیشنهادی این فصل یک جواب 5 -تقریب برای مسئله k -مرکز وزن دار با ظرفیت‌های نرم می‌باشد.

^{۱۰} این موضوع از آن جهت اهمیت دارد که اگر اندازه دامنه یک پادشاه کاهش پیدا کند، ممکن است یک مرکز به دو مرکز تبدیل شود که چنین چیزی در الگوریتم ما وجود ندارد

فصل ۵

الگوریتم‌های تقریبی جدید برای مسئله‌ی k -مرکز ظرفیت‌دار در مدل توزیع شده

در این فصل سایر نتایج جدید پژوهش را، که به‌طور ویژه بر حل مسئله‌ی k -مرکز با ظرفیت‌های نرم و سخت (تعریف ۱-۲) و در مدل توزیع شده تمرکز دارند ارائه می‌دهیم. در تمام الگوریتم‌ها از شیوه مجموعه‌های هسته‌ی ترکیب‌پذیر استفاده شده است. بدین منظور در بخش نخست به معرفی مفاهیم اولیه می‌پردازیم. سپس یک الگوریتم تقریبی توزیع شده برای حل مسئله‌ی k -مرکز با ظرفیت‌های نرم در مدل توزیع شده ارائه می‌دهیم. سپس با فرض وجود یک الگوریتم برای حل مسئله‌ی k -مرکز وزن‌دار با ظرفیت‌های نرم، یک نسخه‌ی بهبودیافته از الگوریتم پیشین خود ارائه می‌دهیم. در ادامه نحوه استفاده از الگوریتم ارائه شده را برای حل مسئله‌ی k -مرکز با ظرفیت‌های سخت بررسی می‌کنیم. در انتها، میزان بهبود به‌دست آمده را به کمک مقایسه نتایج به‌دست آمده با کارهای پیشین نمایش می‌دهیم.

۱-۵ تعاریف و مفاهیم اولیه

برای راحتی در بخش‌های بعدی، چند تعریف و ساختار نوشتاری را در اینجا معرفی کنیم.

- اگر مجموعه‌ی مرجع U ، فضای کل نقاط ممکن تصور شود، برای هر دو نقطه‌ی $a, b \in U$ ،

فاصله‌ی a و b را با $d(a, b)$ نمایش می‌دهیم:

$$d : U \times U \rightarrow \mathbb{R}^+$$

• در این گزارش فرض بر این است که فاصله‌ی نقاط دارای خاصیت متریک است.

$$\forall a, b, c \in U : d(a, c) \leq d(a, b) + d(b, c) \quad (۱ - ۵) \quad \text{یعنی:}$$

$$\forall a, b \in U : d(a, b) = d(b, a) \quad \text{و}$$

• برای نقطه‌ی $a \in U$ و مجموعه نقاط $B \subseteq U$:

$$d(a, B) = \min_{b \in B} \{d(a, b)\} \quad \text{منظور از } d(a, B), \text{ فاصله‌ی } a \text{ تا نزدیک‌ترین عضو } B \text{ است؛ یعنی:}$$

$$\forall b \in B : d(a, B) \leq d(a, b) \quad \text{پس:}$$

تعریف ۱-۵ منظور از $\delta(a, B)$ ، نزدیک‌ترین عضو B به a است.

$$\delta(a, B) = \arg \min_{b \in B} \{d(a, b)\} \quad \text{یعنی:}$$

$$d(a, B) = d(a, \delta(a, B)) \quad \text{پس:}$$

$$\forall b \in B : d(a, \delta(a, B)) \leq d(a, b) \quad \text{و}$$

تعریف ۲-۵ مجموعه نقاط منتسب شده به یک مرکز مانند c را با \mathcal{D}_c نمایش می‌دهیم و فرض می‌کنیم که $c \in \mathcal{D}_c$ است.

تعریف ۳-۵ برای نقطه دلخواه p ، منظور از $\phi(p)$ مرکزی است که p به آن منتسب شده است. در نتیجه

$$p \in \mathcal{D}_c \iff \phi(p) = c \quad \text{به ازای هر مرکز مانند } c, \text{ داریم } c \in \phi(c). \text{ همچنین داریم:}$$

تعریف ۴-۵ با توجه به تعریف مسئله « k -مرکز با ظرفیت‌های نرم» (تعریف ۱-۲)، یک مجموعه مانند

$$D = \{d_1, d_2, \dots, d_m\}$$

که شرایط زیر وجود داشته باشد:

$$1. m \leq k$$

$$2. \forall p \in S : \exists d_i \in D \mid p \in \mathcal{D}_{d_i}$$

$$\forall d_i \in D : |D_{d_i}| \leq L \quad ۳$$

تعریف ۵-۵ منظور از r_S^* شعاع بهینه مسئله « k -مرکز با ظرفیت‌های نرم» به ازای ورودی S است.

۲-۵ الگوریتم تقریبی برای مسئله k -مرکز با ظرفیت‌های نرم

در این بخش یک الگوریتم تقریبی برای مسئله k -مرکز با ظرفیت‌های نرم در مدل توزیع شده ارائه خواهیم داد. ایده به کاررفته در این الگوریتم، استفاده از مجموعه‌های هسته‌ی ترکیب‌پذیر برای ایجاد قابلیت توزیع شدگی در یک الگوریتم موجود و غیر توزیع شده از مسئله k -مرکز با ظرفیت‌های نرم است. خروجی الگوریتم ارائه شده در این بخش یک مجموعه‌ی هسته شامل k نقطه خواهد بود و مراکز نهایی و نحوه انتساب سایر رئوس به این مراکز مشخص نمی‌گردد. اثبات خواهیم کرد که k مرکز در این مجموعه‌ی هسته وجود دارد که یک جواب ۹-تقریب از جواب بهینه است.

ما در بخش ۳-۵ یک الگوریتم کامل‌تر ارائه خواهیم کرد که نه تنها ضریب تقریب بهتری دارد، بلکه علاوه بر k مرکز نهایی، نحوه انتساب سایر نقاط به این مراکز را نیز مشخص می‌کند؛ اما اگر فرض کنیم که هدف نهایی، تنها پیدا کردن یک مجموعه‌ی هسته k نقطه‌ای از مراکز باشد، این الگوریتم یک الگوریتم ساده و کارا برای این هدف خواهد بود.

پیش از ارائه الگوریتم و معرفی قضیه مربوط به آن، دو لم بسیار مهم و کلیدی که در ادامه مورد استفاده خواهد بود را معرفی می‌کنیم.

۱-۲-۵ رابطه‌ی جواب مسئله k -مرکز و مسئله k -مرکز با ظرفیت‌های نرم

لم ۱-۵ هر جواب α -تقریب برای مسئله k -مرکز، یک مجموعه‌ی هسته‌ی $(\alpha + 1)$ -تقریب برای جواب مسئله k -مرکز با ظرفیت‌های نرم است.

اثبات. در اثبات این بخش از شیوه‌ای مشابه آنچه در کار ضرابی‌زاده و سایرین [۱۰] دیده می‌شود استفاده می‌کنیم و از مقایسه نقاط موجود در مجموعه جواب دو مسئله بهره می‌بریم. مجموعه $C = \{q_1, q_2, \dots, q_k\}$ را بعنوان یک جواب α -تقریب برای مسئله k -مرکز به ازای ورودی S

در نظر بگیرید و شعاع جواب بهینه را r^* بنامیم. همچنین نماد $\phi(p)$ را مرکزی که نقطه‌ی $p \in S$ در جواب C به آن متصل شده است در نظر می‌گیریم.

بطور مشابه، مجموعه $C_{cap}^* = \{q_1, q_2, \dots, q_k\}$ که $C_{cap}^* \subseteq S$ را بعنوان جواب بهینه‌ی مسئله‌ی k -مرکز با ظرفیت‌های نرم به شعاع r_{cap}^* در نظر بگیرید و نماد $\phi_{cap}(p)$ را مرکزی که نقطه‌ی $p \in S$ در جواب C_{cap}^* به آن متصل شده است در نظر می‌گیریم.

به ازای هر مرکز $o \in C_{cap}^*$ یک مرکز $o' = \phi(o)$ در نظر بگیرید و مجموعه مراکز جدید را O بنامید. واضح است که O یک زیرمجموعه برای C است. اگر به ازای هر مرکز $o \in C_{cap}^*$ تمام نقاط متصل شده به o را به $\phi(o)$ (مرکز جدید) منتقل کنیم، O یک جواب جدید برای مسئله‌ی k -مرکز با ظرفیت‌های نرم به ازای ورودی S خواهد بود. همچنین داریم:

$$\begin{aligned} \forall p \in S : d(p, \phi(\phi_{cap}(p))) &\leq r_{cap}^* + \alpha \cdot r^* && \text{(نامساوی مثلثاتی)} \\ &\leq r_{cap}^* + \alpha \cdot r_{cap}^* && (r^* \leq r_{cap}^*) \\ &\leq (1 + \alpha) \cdot r_{cap}^* \end{aligned}$$

در این نگاشت ممکن است چند نقطه از C_{cap}^* به یک نقطه یکسان از C نگاشت شوند. ولی از آنجاکه در این نسخه از مسئله ظرفیت‌ها نرم است، این نحوه نگاشت مشکلی ایجاد نمی‌کند و طبق تعریف ۴-۵، مجموعه O یک جواب قابل قبول برای مسئله‌ی k -مرکز با ظرفیت‌های نرم است. \square

۲-۲-۵ یک لم کلیدی

لم ۲-۵ اگر r_S^* شعاع جواب بهینه مسئله‌ی k -مرکز با ظرفیت‌های نرم به ازای مجموعه نقاط ورودی S باشد، به ازای هر مجموعه $\bar{S} \subseteq S$ داریم:

$$r_{\bar{S}}^* \leq 2 \cdot r_S^*$$

اثبات. از آنجاکه در این قسمت دو جواب بهینه (یکی به ازای مجموعه S و دیگری به ازای مجموعه \bar{S}) مورد بحث است، تعریف ۲-۵ و ۳-۵ را اندکی گسترش می‌دهیم.

فرض می‌کنیم به ازای هر نقطه $p \in \bar{S}$: نماد $\phi(p)$ مرکزی که p در جواب بهینه‌ی S به آن منتسب شده است و نماد $\bar{\phi}(p)$ مرکزی که p در جواب بهینه‌ی \bar{S} به آن منتسب شده است را نشان می‌دهد. همچنین فرض می‌کنیم که به ازای هر نقطه (مرکز) مانند c : نماد D_c مجموعه نقاطی که در جواب بهینه‌ی S به c منتسب شده‌اند و نماد \bar{D}_c مجموعه نقاطی که در جواب بهینه‌ی \bar{S} به c منتسب شده‌اند را نشان می‌دهد.

پس اگر مجموعه $C^* = \{q_1, q_2, \dots, q_k\}$ که $C^* \subseteq S$ جواب بهینه برای ورودی S باشد، می‌دانیم:

$$S = \bigcup_{q_i \in C^*} D_{q_i}$$

قصد داریم یک مجموعه مانند \bar{C} معرفی نماییم که یک جواب قابل قبول (و نه بهینه) برای ورودی \bar{S} باشد. بدین منظور تابع $f: C^* \rightarrow \bar{S}$ را به صورت زیر تعریف می‌کنیم:

$$\forall i \in \{1, 2, \dots, k\}: f(q_i) = \delta(q_i, \bar{S}) \quad (5-2)$$

حال بر اساس تابع f ، یک جواب قابل قبول را با تعریف سه عنصر مجموعه $\bar{C} = \{\bar{q}_1, \bar{q}_2, \dots, \bar{q}_k\}$ (که $\bar{C} \subseteq \bar{S}$) و \bar{D} و $\bar{\phi}$ به صورت زیر تعریف می‌کنیم:

$$\forall \bar{q}_i \in \bar{C}: \bar{q}_i = f(q_i)$$

$$\forall p \in \bar{S}: \bar{\phi}(p) = f(\phi(p))$$

$$\forall \bar{q}_i \in \bar{C}: \bar{D}_{\bar{q}_i} = \{p \in \bar{S} | p \in D_{f^{-1}(\bar{q}_i)}\}$$

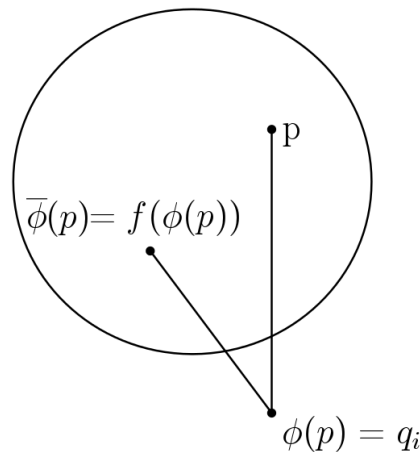
در این نگاشت ممکن است چند نقطه از C^* به یک نقطه یکسان از \bar{S} نگاشت شوند. ولی از آنجا که در این نسخه از مسئله ظرفیت‌ها نرم است، این نحوه نگاشت مشکلی ایجاد نمی‌کند و طبق تعریف ۴-۵، مجموعه \bar{C} یک جواب قابل قبول برای \bar{S} با شعاع $r_{\bar{S}} = \max_{p \in \bar{S}} d(p, \bar{\phi}(p))$ است. همچنین طبق تعریف ۱-۵ و رابطه ۲-۵ می‌دانیم:

$$\forall p \in \bar{S}: d(\phi(p), \bar{\phi}(p)) = d(\phi(p), \delta(\phi(p), \bar{S})) \leq d(\phi(p), p) \leq r_S^* \quad (5-3)$$

مطابق شکل ۱-۵، به ازای تمام نقاط $p \in \bar{S}$ داریم:

$$\forall p \in \bar{S}: d(p, \bar{\phi}(p)) \leq d(p, \phi(p)) + d(\phi(p), \bar{\phi}(p)) \quad (\text{طبق } 1-5)$$

$$\forall p \in \bar{S}: d(p, \bar{\phi}(p)) \leq 2 \cdot r_S^* \quad (\text{طبق } 3-5)$$



شکل ۵-۱: رابطه هندسی بین سه نقطه $p \in \bar{S}$ و $\phi(p)$ و $\bar{\phi}(p)$

که به معنی آن است که درواقع فاصله هر نقطه از مرکز در جواب \bar{C} می‌تواند حداکثر $2 \cdot r_S^*$ است. به عبارت دیگر:

$$r_{\bar{S}} \leq 2 \cdot r_S^* \quad (5-4)$$

از طرفی می‌دانیم که به ازای هر جوابی مانند $r_{\bar{S}}$ داریم:

$$r_S^* \leq r_{\bar{S}} \quad (5-5)$$

از دو رابطه ۵-۴ و ۵-۵ می‌توان نتیجه گرفت:

$$\square \quad r_S^* \leq 2 \cdot r_S^*$$

۳-۲-۵ الگوریتم پیشنهادی

روش پیشنهادی ما برای حل مسئله‌ی k -مرکز با ظرفیت‌های نرم به شرح زیر است:

۱. ابتدا مجموعه نقاط ورودی مسئله (S) را به شکل دلخواه به زیرمجموعه‌های S_1 تا S_m افراز کن و هر افراز را به یک ماشین بده. مجموعه نقاط اختصاص داده شده به ماشین i ام را S_i می‌نامیم.

۲. سپس الگوریتم گنزالس را برای حل مسئله‌ی k -مرکز را بر روی هر ماشین اجرا کن و مجموعه‌ی هسته حاصل را T_i بنام. می‌دانیم که طبق تعریف، مجموعه T_i یک جواب ۲-تقریب برای مجموعه S_i بشمار می‌آید. همچنین به ازای هر $v \in S_i$ ، $\phi_{S_i}(v)$ را برابر با رأسی از T_i در نظر می‌گیریم که v به آن متصل شده است.

۳. بار دیگر الگوریتم گنزالس را بر روی مجموعه T (اجتماع مجموعه‌ها T_i) اجرا کن و مجموعه حاصل را Q بنام. همچنین به ازای هر $v \in T$ ، $\phi_T(v)$ را برابر با رأسی از Q در نظر می‌گیریم که v به آن متصل شده است.

۴. مجموعه Q یک مجموعه‌ی هسته‌ی شامل جواب ۹-تقریب است.

برای اثبات درستی الگوریتم فوق، نیاز است تا اثبات کنیم که k مرکز بر روی نقاط موجود در مجموعه Q وجود دارد که یک جواب ۹-تقریب برای مجموعه نقاط ورودی S است. برای اثبات قضیه ابتدا اثبات می‌کنیم که مجموعه Q یک جواب ۸-تقریب برای مسئله‌ی k -مرکز است، و در انتها از لم ۵-۱ برای اثبات وجود یک جواب ۹-تقریب استفاده می‌کنیم.

قضیه ۳-۵ خروجی الگوریتم فوق شامل یک جواب ۹-تقریب برای مسئله‌ی k -مرکز با ظرفیت‌های نرم است.

اثبات. ابتدا اثبات می‌کنیم که مجموعه‌ی Q یک جواب ۸-تقریب برای مسئله‌ی k -مرکز به ازای ورودی S است. با توجه به گام دوم الگوریتم، می‌دانیم که مجموعه T_i مجموعه‌ی حاصل از اجرای الگوریتم ۲-تقریب گزالس بر روی بخش S_i از ورودی است که توسط ماشین i ام صورت می‌پذیرد. همان‌طور که در بدنه الگوریتم نیز بیان شد، نماد $\phi_{S_i}(p)$ را مشابه با تعریف اصلی (۳-۵)، مرکزی که نقطه‌ی $p \in S_i$ در جواب T_i به آن متصل شده است در نظر می‌گیریم. همچنین شعاع بهینه به ازای بخش S_i را $r_{S_i}^*$ می‌نامیم. با توجه به تعاریف فوق و الگوریتم استفاده‌شده داریم:

$$\forall p \in S_i : d(p, \phi_{S_i}(p)) \leq 2 \cdot r_{S_i}^*$$

همچنین از طرفی به کمک لم ۵-۲ می‌دانیم:

پس داریم:

$$\forall p \in S_i : d(p, \phi_{S_i}(p)) \leq 4 \cdot r_S^* \quad (5-6)$$

به‌طور مشابه و با توجه به گام سوم الگوریتم نیز خواهیم داشت:

$$\forall t \in T : d(t, \phi_T(t)) \leq 4 \cdot r_S^* \quad (5-7)$$

پس:

$$\forall p \in S_i : d(p, \phi_T(\phi_{S_i}(p))) \leq 4 \cdot r_S^* + 4 \cdot r_S^* = 8 \cdot r_S^*$$

و در نتیجه می‌توان گفت که به با توجه به توابع ϕ_{S_i} و ϕ_T ، به ازای هر نقطه از ورودی، یک مرکز در مجموعه Q وجود دارد که فاصله آن حداکثر $8 \cdot r_S^*$ است. حال با توجه به اینکه مجموعه Q یک جواب ۸-تقریب برای مسئله‌ی k -مرکز به ازای ورودی S است و همچنین با توجه به لم ۵-۱ می‌توان نتیجه

گرفت که مجموعه‌ی Q یک مجموعه‌ی هسته ۹-تقریب برای مسئله‌ی k -مرکز با ظرفیت‌های نرم است و k مرکز را می‌توان طوری بر روی نقاط Q قرار داد که یک جواب ۹-تقریب برای مسئله‌ی k -مرکز با ظرفیت‌های نرم باشد. \square

۳-۵ الگوریتم بهبودیافته برای مسئله‌ی k -مرکز با ظرفیت‌های نرم

الگوریتم ارائه شده در بخش ۲-۵ یک الگوریتم ۹-تقریب برای مسئله‌ی k -مرکز با ظرفیت‌های نرم است. در بعضی کاربردها لازم است تا علاوه بر k مرکز نهایی، نحوه انتساب مجموعه نقاط ورودی به این مراکز نیز مشخص شود. به همین منظور در این بخش یک الگوریتم جدید معرفی خواهد شد که علاوه ارائه k مرکز نهایی، نحوه انتساب رئوس را نیز مشخص می‌کند.

در پیاده‌سازی الگوریتم پیشنهادی این بخش از یک الگوریتم β -تقریب برای حل مسئله‌ی k -مرکز وزن‌دار با ظرفیت‌های نرم استفاده می‌شود و یک جواب $2\beta + 4$ ارائه خواهد داد که نسبت به الگوریتم ارائه شده در بخش ۲-۵ بهتر است. مراحل الگوریتم پیشنهادی به صورت زیر است:

۱. ابتدا مجموعه نقاط ورودی مسئله (S) را به شکل دلخواه به زیرمجموعه‌های S_1 تا S_m افراز کن و هر افراز را به یک ماشین بده. مجموعه نقاط اختصاص داده شده به ماشین i ام را S_i می‌نامیم.

۲. سپس الگوریتم ۲-تقریب گنزالس را بر روی هر ماشین اجرا کن و مجموعه‌ی هسته حاصل را T_i بنام و به هر رأس $t \in T_i$ وزنی معادل تعداد نقاط منتسب شده به آن مرکز در نظر بگیر. می‌دانیم که طبق تعریف، مجموعه T_i یک جواب ۲-تقریب برای مجموعه S_i بشمار می‌آید. همچنین به ازای هر $v \in S_i$ ، $\phi_{S_i}(v)$ را برابر با رأسی از T_i در نظر می‌گیریم که v به آن متصل شده است.

۳. یک الگوریتم آفلاین β -تقریب برای حل مسئله‌ی k -مرکز وزن‌دار با ظرفیت‌های نرم (مانند الگوریتم ارائه شده در بخش ۲-۴) بر روی مجموعه وزن‌دار T اجرا کن و مجموعه حاصل را Q بنام. همچنین به ازای هر $v \in T$ ، $\phi_T(v)$ را برابر با رأسی از Q در نظر می‌گیریم که v به آن متصل شده است.

۴. مجموعه Q یک جواب $(2\beta + 4)$ -تقریب از جواب بهینه خواهند بود. همچنین به ازای هر رأس $v \in S_i$ ، رأس v را به مرکز Q $\phi_T(\phi_{S_i}(v)) \in Q$ متصل کن.

برای اثبات درستی الگوریتم فوق، ما نیاز داریم تا اثبات نماییم که k نقطه‌ی موجود در مجموعه Q یک جواب $(4 + 2\beta)$ -تقریب برای مجموعه نقاط ورودی S است. این ادعا در قضیه ۴-۵ آورده شده است که در ادامه بیان می‌کنیم.

قضیه ۴-۵ به ازای ورودی S برای مسئله‌ی k -مرکز با ظرفیت‌های نرم، خروجی الگوریتم ۳-۵ یک جواب $(4 + 2\beta)$ -تقریب از جواب بهینه است.

اثبات. برخلاف اثبات ارائه شده برای قضیه ۳-۵، در اثبات این قضیه از شیوه مقایسه خروجی الگوریتم با نقاط موجود در جواب بهینه استفاده نخواهیم کرد و اثباتی مستقیم ارائه خواهیم داد. با توجه به یکسان بودن گام دوم الگوریتم با الگوریتم قبلی، از تکرار بخشی از اثبات، تعاریف و روابط مربوط به این مرحله از الگوریتم خودداری می‌کنیم. به همین منظور تعریف تابع ϕ_{S_i} و g و شعاع $r_{S_i}^*$ و رابطه ۶-۵ را از اثبات قضیه ۴-۵ صحیح فرض می‌کنیم.

فرض کنیم که الگوریتم بکار گرفته شده برای حل مسئله‌ی k -مرکز وزن‌دار با ظرفیت‌های نرم در گام سوم از الگوریتم پیشنهادی را ALG بنامیم. اگر فرض کنیم که الگوریتم ALG به درستی عمل می‌کند، آنگاه جواب الگوریتم پیشنهادی نیز مشخصاً یک جواب صحیح برای مسئله خواهد بود. این استدلال از آن جهت درست است که ما می‌توانیم هر نقطه از مجموعه وزن‌دار T را نماینده تعدادی نقطه دیگر فرض کنیم و اگر الگوریتم ALG بتواند یک نقطه مانند $p \in T$ (با وزن r_p) را به کمک تعدادی مرکز پوشش دهد، می‌توان نقاطی که در دامنه نقطه p (D_p به دست آمده از گام اول الگوریتم) قرار دارند (و حداکثر تعدادشان r_p است) را به این مراکز $\phi(p)$ منتسب کرد. همچنین از آنجا که فرض می‌شود که الگوریتم ALG به هر عضو Q حداکثر L واحد از رئوس وزن‌دار مجموعه T منتسب کرده است، امکان نقض ظرفیت‌ها وجود ندارد.

اگر به ازای هر رأس $p \in S$ فرض کنیم که $t = \phi_{g(p)}(p)$ (که $t \in T$) آنگاه در رابطه با شعاع جواب خروجی الگوریتم فوق داریم:

$$\begin{aligned}
 d(p, \phi_T(t)) &\leq d(p, t) + d(t, \phi_T(t)) \\
 &\leq 2.2.r_S^* + d(t, \phi_T(t)) && (\text{طبق ۵-۶}) \\
 &\leq 4.r_S^* + \beta.r_T^* && (\text{طبق تعریف } ALG) \\
 &\leq 4.r_S^* + 2.\beta.r_S^* && (\text{طبق قضیه ۵-۲}) \\
 &= (4 + 2\beta).r_S^*
 \end{aligned}$$

پس مجموعه Q یک جواب $(4 + 2\beta)$ -تقریب از جواب بهینه است.

□

نتیجه‌ی ۵-۵ اگر از الگوریتم ۵-تقریب ارائه شده در بخش ۴-۲ در پیاده‌سازی الگوریتم پیشنهادی فوق استفاده شود، یک الگوریتم ۱۴-تقریب توزیع‌شده برای مسئله‌ی k -مرکز با ظرفیت‌های نرم به دست می‌آید.

۴-۵ الگوریتم تقریبی برای مسئله‌ی k -مرکز با ظرفیت‌های سخت

از آنجاکه در بسیاری از کاربردها امکان «نرم» در نظر گرفتن ظرفیت‌ها وجود ندارد، در این بخش قصد داریم تا یک الگوریتم جدید بر مبنای الگوریتم بخش ۵-۳ برای حل مسئله k -مرکز با ظرفیت‌های سخت ارائه دهیم. در واقع در این بخش ما یک شیوه کلی ارائه خواهیم نمود که به کمک آن می‌توان هر جواب مسئله‌ی k -مرکز با ظرفیت‌های نرم را به یک جواب برای مسئله‌ی k -مرکز با ظرفیت‌های سخت تبدیل نمود. این شیوه را قالب قضیه‌ی زیر بیان می‌کنیم.

قضیه‌ی ۵-۶ هر جواب α -تقریب برای مسئله‌ی k -مرکز با ظرفیت‌های نرم را می‌توان به یک جواب (2α) -تقریب برای مسئله‌ی k -مرکز با ظرفیت‌های سخت تبدیل نمود.

اثبات. فرض کنید مجموعه $C = \{q_1, q_2, \dots, q_k\}$ که $C \subseteq S$ و تابع انتساب $f : S \rightarrow C$ یک جواب بهینه α -تقریب برای نسخه ظرفیت نرم مسئله باشد. لازم به ذکر نیست که ممکن است دو مرکز مانند

$q_i, q_j \in C (i \neq j)$ وجود داشته باشند که $q_i = q_j$. بسته به آنکه یک نقطه یک‌بار یا بیش از یک‌بار به‌عنوان مرکز انتخاب‌شده است، می‌توان مجموعه نقاط خروجی این جواب را به دودسته تقسیم‌بندی نمود:

۱. دسته اول آن نقاطی هستند که تنها یک‌بار به‌عنوان مرکز انتخاب‌شده‌اند. در این صورت حداکثر L نقطه وجود دارند که به چنین مرکزی متصل شده‌اند.

۲. دسته دوم آن نقاطی هستند که بیش از یک‌بار به‌عنوان مرکز انتخاب‌شده‌اند و بیش از L نقطه به این نقطه متصل شده است.

اگر تمام اعضای مجموعه C در دسته اول قرار بگیرند، جواب فعلی یک جواب درست برای نسخه سخت از مسئله نیز هست و هیچ عملی لازم نیست.

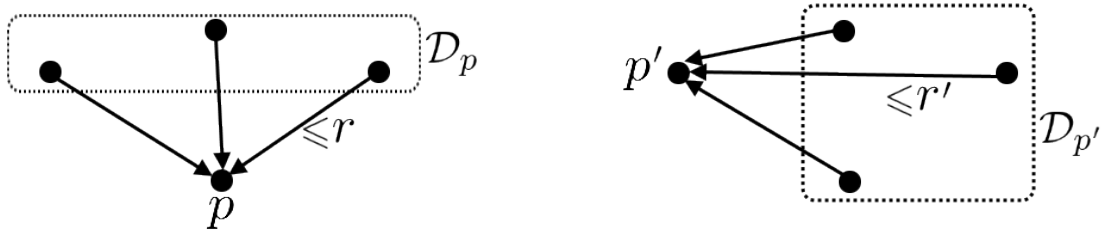
حال فرض کنیم یک نقطه مانند $p \in S$ وجود داشته باشد که t بار در مجموعه C ظاهر شده باشد. برای تمایز بین مراکز واقع‌شده بر روی p ، آن‌ها را مجموعه $C_p = \{p_1, p_2, \dots, p_t\}$ می‌نامیم. در این صورت میدانیم که حداکثر $t.L$ نقطه از مجموعه ورودی در دامنه مراکز C_p قرار دارند. از آنجاکه هدف ما قرار نداشتن هیچ دو مرکزی از جواب بر روی یک نقطه یکسان است، با تکرار گام‌های زیر به ازای هر مرکز $p_i (1 \leq i \leq t)$ ، آن مرکز را با یک نقطه جدید جایگزین می‌کنیم:

۱. دامنه مرکز موردبحث را D_{p_i} می‌نامیم. یک نقطه مانند $p'_i \in D_{p_i}$ که $p'_i \neq p_i$ را در نظر می‌گیریم. میدانیم که چنین نقطه‌ای در مجموعه C قرار نداشته است.

۲. نقطه p'_i را به‌عنوان جایگزین مرکز p_i معرفی می‌کنیم.

۳. مجموعه نقاط D_{p_i} و نقطه p_i را به مرکز جدید متصل می‌کنیم.

از آنجاکه هر نقطه مانند $p \in S$ فقط و فقط عضو دامنه یک مرکز بوده است، هیچ دو مرکز جدید نیز نمی‌توانند یکسان باشد. در نتیجه، با تکرار عمل بالا به ازای تمام مراکز، به مجموعه‌ای از مراکز جدید خواهیم رسید که از یکدیگر کاملاً متمایز هستند و می‌تواند به‌عنوان یک جواب صحیح برای مسئله‌ی k -مرکز با ظرفیت‌های سخت تلقی گردد. شکل ۵-۲ نحوه اجرای گام‌های فوق را به ازای یک مرکز فرضی p با سه رأس در دامنه خود نمایش می‌دهد. در این مثال، مرکز p با یک نقطه از دامنه D_p جایگزین شده است. شکل سمت چپ حالت قبل از جابجایی و شکل سمت راست حالت پس از جابجایی را



شکل ۵-۲: نحوه جایگزینی مرکز برای تبدیل جواب نسخه ظرفیت نرم به سخت

نشان می‌دهد.

همچنین به راحتی می‌توان اثبات کرد که شعاع جواب جدید، حداکثر دو برابر شعاع اولیه خواهد بود. برای اثبات، فرض کنیم که شعاع جواب به ازای مراکز C را r بنامیم. همچنین مطابق با شکل ۵-۲ فرض کنید $p \in C$ یک مرکز قدیمی باشد و ما قصد داریم $p' \in D_p$ را به عنوان مرکز جایگزین انتخاب کنیم. در این صورت بر اساس اصل نامساوی مثلثاتی داریم:

$$\forall v \in D_p : d(v, p') \leq d(v, p) + d(p, p') \leq r + r = 2r$$

□ پس واضح است که فاصله هر نقطه تا مرکز جدیدش حداکثر $2r$ خواهد بود.

نتیجه‌ی ۵-۷ با توجه به قضیه ۵-۶ و نتیجه ۵-۵ می‌توانیم یک جواب ۲۸-تقریب برای مسئله‌ی k -مرکز با ظرفیت‌های سخت به دست بیاوریم.

فصل ۶

نتیجه گیری

در این بخش به جمع بندی نتایج نظری به دست آمده از پژوهش می پردازیم. همچنین در انتها به مقایسه این نتایج با با کارهای موجود می پردازیم. سپس در انتها مجموعه پیشنهادهایی برای کارهای آتی ارائه می شود.

۱-۶ نتایج بدست آمده

با فرض وجود یک الگوریتم β -تقریب برای مسئله k -مرکز وزن دار با ظرفیت های نرم، موارد زیر نتایج به دست آمده در این پژوهش را خلاصه می کند:

۱. یک الگوریتم ۵-تقریب برای مسئله k -مرکز وزن دار با ظرفیت های نرم (قضیه ۴-۵)
۲. یک الگوریتم ۹-تقریب برای مسئله k -مرکز با ظرفیت های نرم در مدل توزیع شده و به کمک مجموعه های هسته ی ترکیب پذیر (قضیه ۵-۳).
۳. یک الگوریتم $(4 + 2\beta)$ -تقریب برای مسئله k -مرکز با ظرفیت های نرم در مدل توزیع شده و به کمک مجموعه های هسته ی ترکیب پذیر (قضیه ۵-۴)
۴. یک الگوریتم $(8 + 4\beta)$ -تقریب برای مسئله k -مرکز با ظرفیت های سخت در مدل توزیع شده و به کمک مجموعه های هسته ی ترکیب پذیر (نتیجه ۵-۷)

با توجه به قضیه ۴-۶ داریم $\beta = 5$. در نتیجه به ضریب تقریب‌های جدول ۶-۱ برای نسخه‌های مختلف مسئله k -مرکز خواهیم رسید. لازم به ذکر است که تمام این الگوریتم‌ها برای مدل توزیع شده طراحی شده‌اند.

مسئله	ضریب تقریب	مشخص شدن نحوه انتساب نقاط به مراکز
k -مرکز با ظرفیت‌های نرم	۹	خیر
k -مرکز با ظرفیت‌های نرم	۱۴	بله
k -مرکز با ظرفیت‌های سخت	۲۸	بله

جدول ۶-۱: ضریب تقریب‌های بدست آمده برای نسخه‌های مختلف مسئله k -مرکز ظرفیت‌دار

همچنین جدول ۶-۲ مقایسه‌ای بین الگوریتم ارائه شده در این پژوهش برای حل مسئله k -مرکز با ظرفیت‌های سخت در مدل توزیع شده (بخش ۵-۴) و بهترین کار موجود (پژوهش [۱۲]) ارائه می‌دهد. مشخص است که الگوریتم پیشنهادی در این پژوهش به‌طور مناسبی ضریب تقریب بهترین الگوریتم ماقبل خود را بهبود داده است.

همچنین از نقطه نظر تعداد دور لازم (در مدل نگاشت-کاهش)، الگوریتم پیشنهادی ما با تنها یک دور، بهترین نتیجه را ارائه می‌دهد. نکته حائز اهمیت دیگر آن است که الگوریتم پیشنهادی ما در این پژوهش یک الگوریتم کاملاً قطعی است، در حالی که الگوریتم [۱۲] تصادفی بوده و ضریب تقریب یا تعداد دورهای آن در بدترین حالت بسیار بیشتر خواهد بود.

الگوریتم	ضریب تقریب	دور در نگاشت-کاهش	قطعی/تصادفی
پژوهش جاری	۲۸	۱	قطعی
الگوریتم پژوهش [۱۲]	۶۴	$O(1)$	تصادفی

جدول ۶-۲: مقایسه الگوریتم پیشنهادی با نمونه موجود برای مسئله k -مرکز ظرفیت‌دار

۲-۶ کارهای آینده

در این پژوهش سعی شد که مسئله‌ی k -مرکز با ظرفیت‌های نرم، k -مرکز با ظرفیت‌های سخت و نسخه وزن‌دار آن مورد مطالعه قرار بگیرد. با توجه به اهمیت پردازش داده‌های حجیم و آنکه هدف اصلی ما در این پژوهش طراحی الگوریتم‌هایی برای مدل توزیع شده بود، در طراحی الگوریتم‌های ارائه شده از روش مجموعه‌های هسته ترکیب پذیر استفاده شد. با وجود مزایای مختلف این روش، استفاده از این روش برای حل مسائل خوشه‌بندی بسیار محدود بوده است و به‌طور خاص در کار اخیر باطنی و همکاران [۱۲] قابل مشاهده است. در نتیجه به کارگیری این شیوه برای ارائه الگوریتم‌های توزیع شده جدید با ضرایب تقریب ثابت برای سایر مسائل خوشه‌بندی از قبیل k -میانه، k -میانگین می‌تواند مورد توجه پژوهشگران قرار بگیرد.

هرچند الگوریتم ۴-۵ برای حل مسئله k -مرکز با ظرفیت‌های سخت در مدل توزیع شده از بهترین الگوریتم ماقبل خود بهتر عمل می‌کند، در طراحی این الگوریتم از یک رویکرد دومرحله‌ای استفاده شده است که امکان بهبود آن وجود دارد. در مرحله اول، مسئله با فرض نرم بودن ظرفیت‌ها حل می‌شود و سپس در مرحله دوم حاصل به جوابی برای نسخه‌ی با ظرفیت‌های سخت تبدیل می‌گردد. طراحی یک الگوریتم مستقیم برای حل مسئله‌ی k -مرکز با ظرفیت‌های سخت در مدل توزیع شده می‌تواند به‌عنوان یکی دیگر از پژوهش‌های آتی در نظر گرفته شود. استفاده از الگوریتم ۶-تقریب خولر و سوسمان [۲] برای مسئله‌ی k -مرکز با ظرفیت‌های سخت می‌تواند به انجام این پژوهش کمک کند.

همچنین در این پژوهش فرض یکسان بودن ظرفیت تمام نقاط از ابتدا در تمام مسائل k -مرکز ظرفیت‌دار نظر گرفته شد. حال آنکه طراحی الگوریتم‌هایی برای نسخه‌های با ظرفیت‌های غیر یکسان مسئله k -مرکز می‌تواند به‌عنوان بخش دیگری از کارهای آتی مورد توجه قرار بگیرد.

کتاب نامه

- [1] M. R. Garey and D. S. Johnson. A guide to the theory of NP-completeness. *WH Freeman, New York*, 1979.
- [2] S. Khuller and Y. J. Sussmann. The capacitated k -center problem. *SIAM Journal on Discrete Mathematics*, 13(3):403–418, 2000.
- [3] B. Chen, R. Dutta, and G. N. Rouskas. On the application of k -center algorithms to hierarchical traffic grooming. In *Proceedings of 2nd International Conference on Broadband Networks, 2005.*, pages 1218–1224. IEEE, 2005.
- [4] T. F. Gonzalez. Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science*, 38:293–306, 1985.
- [5] M. Cygan, M. Hajiaghayi, and S. Khuller. LP rounding for k -centers with non-uniform hard capacities. In *Proceedings of the 53th IEEE Symposium on Foundations of Computer Science*, pages 273–282. IEEE, 2012.
- [6] H.-C. An, A. Bhaskara, C. Chekuri, S. Gupta, V. Madan, and O. Svensson. Centrality of trees for capacitated k -center. *Mathematical Programming*, 154(1-2):29–53, 2015.
- [7] H. Karloff, S. Suri, and S. Vassilvitskii. A model of computation for MapReduce. In *Proceedings of the 21st ACM-SIAM Symposium on Discrete Algorithms*, pages 938–948. Society for Industrial and Applied Mathematics, 2010.
- [8] A. Ene, S. Im, and B. Moseley. Fast clustering using mapreduce. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 681–689. ACM, 2011.

- [9] S. Im and B. Moseley. Brief announcement: Fast and better distributed mapreduce algorithms for k -center clustering. In *Proceedings of the 27th ACM Symposium on Parallel Algorithms and Architectures*, pages 65–67. ACM, 2015.
- [10] S. Aghamolaei, M. Farhadi, and H. Zarrabi-Zadeh. Diversity maximization via composable coresets. In *Proceedings of the 27th Canadian Conference on Computational Geometry*, 2015.
- [11] P. Indyk, S. Mahabadi, M. Mahdian, and V. S. Mirrokni. Composable core-sets for diversity and coverage maximization. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 100–108. ACM, 2014.
- [12] M. Bateni, A. Bhaskara, S. Lattanzi, and V. Mirrokni. Distributed balanced clustering via mapping coresets. In *Advances in Neural Information Processing Systems*, pages 2591–2599, 2014.
- [13] P. K. Agarwal, S. Har-Peled, and K. R. Varadarajan. Approximating extent measures of points. *Journal of the ACM*, 51(4):606–635, 2004.
- [14] R. Xu and D. Wunsch. Survey of clustering algorithms. *IEEE Transactions on neural networks*, 16(3):645–678, 2005.
- [15] J. Barilan, G. Kortsarz, and D. Peleg. How to allocate network centers. *Journal of Algorithms*, 15(3):385–415, 1993.
- [16] M. Cygan and T. Kociumaka. Constant factor approximation for capacitated k -center with outliers. In *LIPICs-Leibniz International Proceedings in Informatics*, volume 25. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2014.

Abstract

Clustering is one of the most well-known fundamental problems in computer science. In this thesis we have focused on a particular version of such problem, called capacitated k -center, and we analyze and survey them in the distributed model, when massive data is given. Our contribution in this research includes a new approximation algorithms with constant approximate factors for such problems in the distributed model, as well as improving the best available algorithm for capacitated k -center problem. Composable coresets as one of the most important techniques in distributed and streaming model is our primary tools in designing these algorithms. This technique gives the opportunity of solving the problem in smaller chunks of data, and giving the result by combining them.

Keywords: Approximation algorithm, Composable coresets, Clustering, Capacitated k -center, Distributed model



Sharif University of Technology

Department of Computer Engineering

M.Sc. Thesis

**Approximation Algorithms for Clustering Points
in the Distributed Model**

By:

Emad Aghajani

Supervisor:

Dr. Hamid Zarrabi-Zadeh

July 2016