



دانشگاه صنعتی شریف
دانشکده‌ی مهندسی کامپیوتر

پایان‌نامه‌ی کارشناسی ارشد
گرایش مهندسی نرم‌افزار

عنوان:

الگوریتم‌های تقریبی برای خوشه‌بندی نقاط در مدل جویبار داده

نگارش:

بهنام حاتمی ورزنه

استاد راهنما:

حمید ضرابی‌زاده

شهریور ۱۳۹۴



به نام خدا
دانشگاه صنعتی شریف
دانشکده‌ی مهندسی کامپیوتر

پایان‌نامه‌ی کارشناسی ارشد

عنوان: الگوریتم‌های تقریبی برای خوشه‌بندی نقاط در مدل جویبار داده
نگارش: بهنام حاتمی ورزنه

کمیته‌ی ممتحنین

استاد راهنما: حمید ضرابی‌زاده
امضاء:

استاد مشاور: حمید بیگی
امضاء:

استاد مدعو: استاد ممتحن
امضاء:

تاریخ:

چکیده

این قسمت باید تکمیل گردد.

کلیدواژه‌ها: خوشه‌بندی، k -مرکز، جویبار داده، الگوریتم تقریبی

فهرست مطالب

۹	۱ مقدمه
۱۰	۱-۱ تعریف مسئله
۱۳	۲-۱ اهمیت موضوع
۱۴	۳-۱ ادبیات موضوع
۱۴	۴-۱ اهداف تحقیق
۱۵	۵-۱ ساختار پایان نامه
۱۷	۲ مفاهیم اولیه
۱۷	۱-۲ مسائل ان پی- سخت
۱۹	۲-۱ پوشش رأسی
۲۰	۲-۲ الگوریتم های تقریبی
۲۲	۲-۲-۱ میزان تقریب پذیری مسائل
۲۲	۳-۲ الگوریتم های جویبار داده
۲۳	۲-۳-۱ مقدمه
۲۴	۲-۳-۲ گونه های مطرح
۲۴	۳-۳-۲ تحلیل الگوریتم های جویبار داده
۲۵	۲-۳-۴ مجموعه هسته

۲۸	۳ کارهای پیشین
۲۸	۳-۱ k -مرکز در حالت ایستا
۳۲	۳-۲ k -مرکز در حالت جویبار داده
۳۹	۳-۳ k -مرکز با داده‌های پرت
۴۵	۴ نتایج جدید
۴۶	۵ نتیجه‌گیری
۴۷	آ مطالب تکمیلی

فهرست شکل‌ها

۱-۱	نمونه‌ای از مسئله‌ی ۲- مرکز	۱۱
۲-۱	نمونه‌ای از مسئله‌ی ۲- مرکز با داده‌های پرت	۱۲
۳-۱	نمونه‌ای از مسئله‌ی ۲- مرکز در حالت پیوسته	۱۳
۱-۳	نمونه‌ای از تخصیص نقاط به ازای مراکز آبی رنگ	۲۹
۲-۳	نمونه‌ای از تبدیل یک گراف ورودی مسئله‌ی پوشش رأسی به یک ورودی مسئله‌ی k -مرکز (در گراف سمت چپ، یال‌های سیاه وزن ۱ و یال‌های آبی، وزن ۲ دارند)	۲۹
۳-۳	نمونه‌ای از حل مسئله‌ی ۳- مرکز با الگوریتم گنزالز	۳۱
۴-۳	نمونه‌ای از توری‌بندی الگوریتم ضربابی زاده (نقاط آبی، مراکز به دست آمده از الگوریتم تقریبی است). پس از توری‌بندی کافی است برای هر کدام از خانه‌های شبکه‌بندی، تنها یکی را به عنوان نماینده در نظر بگیریم	۳۳
۵-۳	نمونه‌ای از اجرای الگوریتم ضربابی زاده بر روی چهار نقطه $P_1 \dots P_4$ که به ترتیب اندیس در جویبار داده فرا می‌رسند و دایره‌های $B_1 \dots B_4$ دایره‌هایی که الگوریتم به ترتیب نگه می‌دارد	۳۵
۶-۳	اثبات لم ۲-۳	۳۷
۷-۳	اثبات لم ۳-۳	۳۹
۸-۳	کاهش قابل توجه شعاع مسئله‌ی ۲- مرکز با حذف تنها دو نقطه	۴۰

فهرست جدول‌ها

۲-۱ نمونه‌هایی از ضرایب تقریب برای مسائل بهینه‌سازی ۲۲

فصل ۱

مقدمه

مسئله‌ی خوشه‌بندی^۱ یکی از مهم‌ترین مسائل داده‌کاوی^۲ به حساب می‌آید. در این مسئله هدف، دسته‌بندی تعدادی جسم به گونه‌ای است که اجسام در یک دسته (خوشه)، نسبت به یکدیگر در برابر دسته‌های دیگر شبیه‌تر باشند (معیارهای متفاوتی برای تشابه تعریف می‌گردد). این مسئله در حوزه‌های مختلفی از علوم کامپیوتر، از جمله داده‌کاوی، جست‌وجوی الگو^۳، پردازش تصویر^۴، بازیابی اطلاعات^۵ و بیوانفورماتیک^۶ مورد استفاده قرار می‌گیرد [۱].

مسئله‌ی خوشه‌بندی، به‌خودی‌خود یک مسئله‌ی الگوریتمی به حساب نمی‌آید. راه‌حل‌های الگوریتمی بسیار زیادی برای خوشه‌بندی تعریف شده است. این الگوریتم‌ها را براساس رویکردهای مختلفی که به مسئله دارند، می‌توان در یکی از چهار دسته‌بندی زیر قرار داد:

- خوشه‌بندی‌های سلسه‌مراتبی^۷

- خوشه‌بندی‌های مرکزگرا^۸

^۱ Clustering

^۲ Data mining

^۳ Pattern recognition

^۴ Image analysis

^۵ Information retrieval

^۶ Bioinformatics

^۷ Hierarchical clustering

^۸ Centroid-based clustering

• خوشه‌بندی‌های مبتنی بر توزیع^۹ نقاط

• خوشه‌بندی‌های مبتنی بر چگالی^{۱۰} نقاط

در عمل هیچ کدام از راه‌حل‌های بالا بر دیگری ارجحیت ندارند و باید راه‌حل مد نظر را متناسب با کاربرد مطرح مورد استفاده قرار داد. به طور مثال استفاده از الگوریتم‌های مرکز‌گرا، برای خوشه‌های غیر محدب به خوبی عمل نمی‌کند. یکی از راه‌حل‌های شناخته شده برای مسئله‌ی خوشه‌بندی، الگوریتم k -مرکز است. در این الگوریتم هدف، پیدا کردن k نقطه به عنوان مرکز دسته‌ها است به طوری که شعاع دسته‌ها تا حد ممکن کمینه شود. در نظریه‌ی گراف، مسئله‌ی k -مرکز متریک^{۱۱} یا مسئله‌ی استقرار تجهیزات متریک^{۱۲} یک مسئله‌ی بهینه‌سازی ترکیبیاتی^{۱۳} است. فرض کنید که n شهر و فاصله‌ی دویه‌دوی آن‌ها، داده شده است. می‌خواهیم k انبار در شهرهای مختلف بسازیم به طوری که بیش‌ترین فاصله‌ی هر شهری از نزدیک‌ترین انبار به خود، کمینه گردد. در حالت نظریه‌ی گراف آن، این بدان معناست که مجموعه‌ای شامل k رأس انتخاب کنیم به طوری که بیش‌ترین فاصله‌ی هر نقطه از نزدیک‌ترین نقطه‌اش داخل مجموعه‌ی k عضوی کمینه گردد. توجه نمایید که فاصله‌ی بین رئوس باید در فضای متریک^{۱۴} باشند و یا به زبان دیگر، یک گراف کامل داشته باشیم که فاصله‌ها در آن در رابطه‌ی مثلثی^{۱۵} صدق می‌کنند. مثالی از مسئله‌ی ۲-مرکز را در شکل ۱-۱ نشان داده شده است.

در این پژوهش، مسئله‌ی k -مرکز با متریک‌های خاص و برای k های کوچک مورد بررسی قرار گرفته است و از جنبه‌های متفاوتی بهبود یافته است. در زیر فصل بعدی، تعریف رسمی^{۱۶} از مسائلی که در این پایان‌نامه مورد بررسی قرار می‌گیرند را تعریف کرده و در مورد هر کدام توضیح مختصری می‌دهیم.

۱-۱ تعریف مسئله

تعریف دقیق‌تر مسئله‌ی k -مرکز در زیر آمده است:

^۹Distribution-based

^{۱۰}Density-based

^{۱۱}Metric

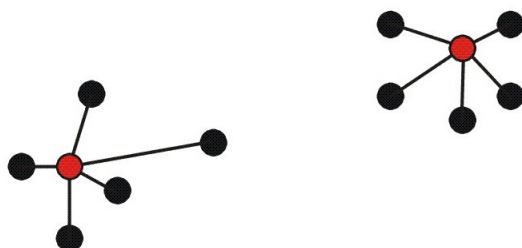
^{۱۲}Metric facility location

^{۱۳}Combinatorial optimization

^{۱۴}Metric space

^{۱۵}Triangle equation

^{۱۶}Formal



شکل ۱-۱: نمونه‌ای از مسئله‌ی ۲-مرکز

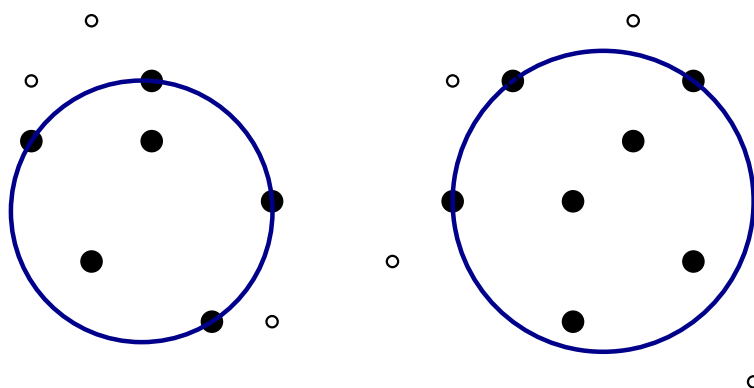
مسئله‌ی ۱-۱ (k -مرکز) یک گراف کامل بدون جهت $G = (V, E)$ با فاصله‌ی d ، که از نامساوی مثلی پیروی می‌کند داده شده است. زیرمجموعه‌ی $S \subseteq V$ با اندازه‌ی k را به گونه‌ای انتخاب کنید به طوری که عبارت زیر را کمینه کند:

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

گونه‌های مختلفی از مسئله‌ی k -مرکز با محدودیت‌های متفاوتی به وسیله‌ی پژوهشگران، مورد مطالعه قرار گرفته است. از جمله‌ی این گونه‌ها، می‌توان به حالتی که در بین داده‌های ورودی، داده‌های پرت وجود دارد، می‌توان اشاره کرد. در واقع در این مسئله، قبل از خوشه‌بندی می‌توانیم تعدادی از نقاط ورودی را حذف نماییم و سپس به خوشه‌بندی نقاط پردازیم. سختی این مسئله از آنجاست که نه تنها باید مسئله‌ی خوشه‌بندی را حل نمود، بلکه در ابتدا باید تصمیم گرفت که کدام یک از داده‌ها را به عنوان داده‌ی پرت در نظر گرفت که بهترین جواب در زمان خوشه‌بندی به دست آید. در واقع اگر تعداد نقاط پرتی که مجاز به حذف است، برابر صفر باشد، مسئله به مسئله‌ی k -مرکز تبدیل می‌شود. نمونه‌ای از مسئله‌ی ۲-مرکز با ۷ داده‌ی پرت را در شکل ۱-۲ می‌توانید ببینید. تعریف دقیق‌تر این مسئله در زیر آمده است:

مسئله‌ی ۱-۲ (k -مرکز با داده‌های پرت) یک گراف کامل بدون جهت $G = (V, E)$ با فاصله‌ی d ، که از نامساوی مثلی پیروی می‌کند داده شده است. زیرمجموعه‌ی $Z \subseteq V$ با اندازه‌ی z و $S \subseteq V - Z$ با اندازه‌ی k را انتخاب کنید به طوری که عبارت زیر را کمینه کند:

$$\max_{v \in (V-Z)} \{ \min_{s \in S} d(v, s) \}$$



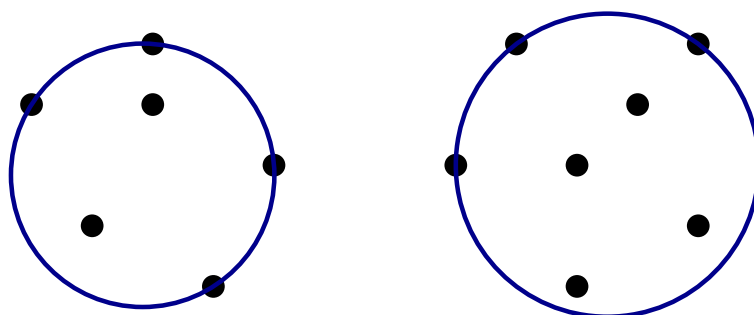
شکل ۱-۲: نمونه‌ای از مسئله‌ی ۲-مرکز با داده‌های پرت

گونه‌ی دیگری از مسئله‌ی k -مرکز که در سال‌های اخیر مورد توجه قرار گرفته است، حالت جویبار داده‌ی آن است. در این گونه از مسئله‌ی k -مرکز، در ابتدا تمام نقاط در دسترس نیستند، بلکه به مرور زمان نقاط در دسترس قرار می‌گیرند. محدودیت دومی که وجود دارد، محدودیت شدید حافظه است، به طوری که نمی‌توان تمام نقاط را در حافظه نگه داشت و بعضاً حتی در حافظه‌ی جانبی نیز نمی‌توان ذخیره نمود و به طور معمول باید مرتبه‌ی حافظه‌ای کم‌تر از مرتبه‌ی حافظه‌ی خطی^{۱۷} متناسب با تعداد نقاط استفاده نمود. از این به بعد به چنین مرتبه‌ای، مرتبه‌ی زیرخطی می‌گوییم. مدلی که ما در این پژوهش بر روی آن تمرکز داریم مدل جویبار داده تک‌گذره [۲] است. یعنی تنها یک بار می‌توان از ابتدا تا انتهای داده‌ها را بررسی کرد و پس از عبور از یک داده، اگر آن را در حافظه ذخیره نکرده باشیم، دیگر به آن دسترسی نداریم.

یکی از دغدغه‌هایی که در مسائل جویبار داده وجود دارد، عدم امکان دسترسی به تمام نقاط است. در واقع هم این مشکل وجود دارد که به تمام داده‌های قبلی دسترسی نداریم و هم این مشکل وجود دارد که هیچ اطلاعی از داده‌های آتی نداریم. در نتیجه یکی از تبعات این گونه از مسئله‌ی k -مرکز، امکان انتخاب نقطه‌ای به عنوان مرکز برای یک دسته به طوری که در بین نقاط ورودی نیست. این گونه از مسئله‌ی k -مرکز، معمولاً تنها برای L_p -متریک مطرح می‌شود یا حالتی که ما مجموعه‌ای از تمام نقاط فضا به انضمام فاصله‌هایشان را داشته باشیم. زیرا مرکز دسته‌ها ممکن است در هر نقطه از فضا قرار بگیرد و ما نیاز داریم که فاصله‌ی آن را از تمام نقاط بدانیم. نمونه‌ای از مسئله‌ی ۲-مرکز در حالت پیوسته، در شکل ۱-۳ نشان داده شده است. تعریف دقیق گونه‌ی جویبار داده‌ی مسئله‌ی k -مرکز، در

^{۱۷}Sublinear

زیر آمده است:



شکل ۱-۳: نمونه‌ای از مسئله‌ی ۲-مرکز در حالت پیوسته

مسئله‌ی ۱-۳ (k -مرکز در حالت جویبار داده) مجموعه‌ی U از نقاط فضای d بعدی داده شده است. زیرمجموعه $S \subseteq U$ با اندازه‌ی k را انتخاب کنید به طوری که عبارت زیر را کمینه کند:

$$\max_{u \in U} \{ \min_{s \in S} L_p(u, s) \}$$

از آنجایی که گونه‌ی جویبار داده و داده پرت مسئله‌ی k -مرکز به علت داغ شدن مبحث داده‌های بزرگ^{۱۸}، به تازگی مورد توجه قرار گرفته است و نتایج به دست آمده قابل بهبود است، در این تحقیق سعی شده است که تمرکز بر روی این گونه‌ی خاص از مسئله باشد. همچنین در این پژوهش سعی می‌شود گونه‌های مسئله را برای انواع متریک‌ها و برای k های کوچک نیز مورد بررسی قرار داد.

۲-۱ اهمیت موضوع

مسئله‌ی k -مرکز و گونه‌های آن کاربردهای زیادی در داده‌کاوی دارند. این الگوریتم یکی از رایج‌ترین الگوریتم‌های مورد استفاده برای خوشه‌بندی محسوب می‌شود. به علت افزایش حجم داده‌ها و تولید داده‌ها در طول زمان، گونه‌ی جویبار داده‌ی مسئله در سال‌های اخیر مورد توجه قرار گرفته است. از طرفی مسئله‌ی k -مرکز در بعضی مسائل مانند ارسال نامه از تعدادی مراکز پستی به گیرنده‌ها، نیاز دارد که تمام نامه‌ها را به دست گیرنده‌ها برساند و در نتیجه باید همه‌ی نقاط را پوشش دهد، ولی برای بعضی از مسائل تجاری، نیازی به پوشش تمام نقاط هدف نیست و از لحاظ اقتصادی به صرفه است که نقاط پرت را در نظر نگرفت. به طور مثال، Kmarts در سال ۲۰۱۰ اعلام کرده است که به ۸۸ درصد جمعیت

^{۱۸} Big data

آمریکا با فاصله‌ای حداکثر ۶ مایل، می‌تواند سرویس دهد. در صورتی که اگر این شرکت قصد داشت تمام جمعیت آمریکا را پوشش دهد، نیاز داشت تعداد شعب یا شعاع پوشش خود را به میزان قابل توجهی افزایش دهد که از لحاظ اقتصادی به صرفه نیست [۳]. علاوه بر دو گونه‌ی مطرح شده در این قسمت، گونه‌های دیگر از مسئله‌ی k -مرکز در مرجع [۳] آمده است.

۳-۱ ادبیات موضوع

همان‌طور که ذکر شد مسئله‌ی k -مرکز در حالت داده‌های پرت و جویبار داده، گونه‌های تعمیم‌یافته از مسئله‌ی k -مرکز هستند و در حالت‌های خاص به مسئله‌ی k -مرکز کاهش پیدا می‌کنند. مسئله‌ی k -مرکز در حوزه‌ی مسائل ان‌پی-سخت^{۱۹} قرار می‌گیرد و با فرض $P \neq NP$ الگوریتم دقیق با زمان چندجمله‌ای برای آن وجود ندارد. بنابراین برای حل کارای^{۲۰} این مسائل از الگوریتم‌های تقریبی^{۲۱} استفاده می‌شود.

برای مسئله‌ی k -مرکز، دو الگوریتم تقریبی معروف وجود دارد. در الگوریتم اول، که به روش حریصانه^{۲۲} عمل می‌کند، در هر مرحله بهترین مرکز ممکن را انتخاب می‌کند به طوری تا حد ممکن از مراکز قبلی دور باشد. این الگوریتم، الگوریتم تقریبی با ضریب تقریب ۲ ارائه می‌دهد. در الگوریتم دوم، با استفاده از مجموعه‌ی غالب کمینه^{۲۳}، الگوریتمی با ضریب تقریب ۲ ارائه می‌گردد. همچنین ثابت شده است، که بهتر از این ضریب تقریب، الگوریتمی نمی‌توان ارائه داد مگر آن‌که $P = NP$ باشد.

۴-۱ اهداف تحقیق

در این پایان‌نامه مسئله‌ی k -مرکز در حالت جویبار داده با داده‌های پرت در حالت‌های مختلف مورد بررسی قرار می‌گیرد و سعی خواهد شد که نتایج قبلی در این مسائل را از جنبه‌های مختلفی مورد بهبود دهد.

^{۱۹} NP-hard

^{۲۰} Efficient

^{۲۱} Approximation Algorithm

^{۲۲} Greedy

^{۲۳} Dominating set

اولین مسئله‌ای که مورد بررسی قرار گرفته است، ارائه الگوریتمی تقریبی، برای مسئله‌ی ۱-مرکز در حالت جویبار داده است به طوری که، نه تنها تمام نقاط ورودی را می‌پوشاند بلکه تضمین می‌کند که دایره‌ی بهینه‌ی جواب ۱-مرکز برای نقاط ورودی را نیز بپوشاند. در تلاش اول، الگوریتم جدیدی با حافظه و زمان به‌روزرسانی $O(\frac{d}{\epsilon})$ و با ضریب تقریب $1/8 + \epsilon$ ، برای این مسئله ارائه گردید. با بررسی‌های بیش‌تر، الگوریتم دیگری با ضریب تقریب $1/69$ ، با الگوریتمی کاملاً متفاوت، با حافظه‌ی $O(d)$ برای این مسئله ارائه گردید.

مسئله‌ی دومی که مورد بررسی قرار گرفت، مسئله‌ی ۱-مرکز در حالت جویبار داده با داده‌های پرت است. در ابتدا الگوریتمی ساده با حافظه‌ی $O(zd)$ و زمان به‌روزرسانی $O(zd \log(z))$ با ضریب تقریب ۲ برای این مسئله ارائه گردید. در بررسی‌های بعدی، با استفاده از نتایج بدست آمده در قسمت قبل، برای z های کوچک، الگوریتمی با حافظه‌ی $O(dz^d)$ با ضریب تقریب $1/69$ برای این مسئله ارائه شد که ضریب تقریب بهترین الگوریتم موجود را، از $1/79$ به $1/69$ کاهش می‌دهد.

مسئله‌ی سومی که مورد بررسی قرار گرفت، مسئله‌ی ۲-مرکز در حالت جویبار داده با داده‌های پرت است. بهترین الگوریتمی که در حال حاضر برای این مسئله وجود دارد، یک الگوریتم با ضریب تقریب $4 + \epsilon$ است [۴]. ما با ارائه الگوریتمی جدید، الگوریتمی با ضریب تقریب $1/8 + \epsilon$ برای این مسئله ارائه دادیم که بهبودی قابل توجه محسوب می‌شود.

۵-۱ ساختار پایان‌نامه

این پایان‌نامه در پنج فصل به شرح زیر ارائه خواهد شد. در فصل دوم به بیان مفاهیم و تعاریف مرتبط با موضوعات مورد بررسی در بخش‌های دیگر خواهیم پرداخت. سعی شده، تا حد امکان با زبان ساده و ارجاع‌های مناسب، پایه‌ی لازم را برای فصول بعدی در این فصل فراهم آوریم. فصل سوم این پایان‌نامه شامل مطالعه و بررسی کارهای پیشین انجام شده مرتبط با موضوع این پایان‌نامه خواهد بود. این فصل در سه بخش تنظیم گردیده است. در بخش اول، مسئله‌ی k -مرکز در حالت ایستا مورد بررسی قرار می‌گیرد. در بخش دوم، حالت جویبار داده‌ی مسئله و مجموعه هسته‌های^{۲۴} مطرح برای این مسئله مورد بررسی قرار می‌گیرد. در نهایت، در بخش سوم، مسئله‌ی k -مرکز با داده‌های پرت مورد بررسی قرار می‌گیرد.

^{۲۴}coreset

در فصل چهارم، نتایج جدیدی که در این پایان نامه به دست آمده است، ارائه می شود. این نتایج شامل الگوریتم های جدید برای مسئله ۱ - مرکز در حالت جویبار داده، مسئله ۱ - مرکز با داده های پرت در حالت جویبار داده و مسئله دو مرکز در حالت جویبار داده با داده های پرت می شود.

در فصل پنجم، به جمع بندی کارهای انجام شده در این پژوهش و ارائه پیشنهادهایی برای انجام کارهای آتی و تعمیم هایی که از راه حل ارائه شده وجود دارد، خواهیم پرداخت.

فصل ۲

مفاهیم اولیه

در این فصل به تعریف و بیان مفاهیم پایه‌ای مورد استفاده در فصل‌های بعد می‌پردازیم. با توجه به مطالب مورد نیاز در فصل‌های آتی، مطالب این فصل به سه بخش، مسائل ان‌پی-سخت، الگوریتم‌های تقریبی و الگوریتم‌های جویبار داده تقسیم می‌شود.

۲-۱ مسائل ان‌پی-سخت

یکی از اولین سوال‌های بنیادی مطرح در علم کامپیوتر، اثبات عدم حل پذیری بعضی از مسائل است. به عنوان نمونه، می‌توان از دهمین مسئله‌ی هیلبرت^۱ در گنگره‌ی ریاضی یاد کرد. هیلبرت این مسئله را این‌گونه بیان کرد: ”فرآیندی طراحی کنید که در تعداد متناهی گام بررسی کند که آیا یک چندجمله‌ای، ریشه‌ی صحیح^۲ دارد یا خیر“. با مدل محاسباتی که به‌وسیله‌ی تورینگ^۳ ارائه شد، این مسئله معادل پیدا کردن الگوریتمی برای این مسئله است که اثبات می‌شود امکان‌پذیر نیست [۵]. برخلاف مثال بالا، عمده‌ی مسائل علوم کامپیوتر از نوع بالا نیستند و برای طیف وسیعی از آن‌ها، الگوریتم‌های پایان‌پذیر وجود دارد. بیشتر تمرکز علوم کامپیوتر هم بر روی چنین مسائلی است.

اگر چه برای اکثر مسائل، الگوریتمی پایان‌پذیر وجود دارد، اما وجود چنین الگوریتمی لزومی بر حل

^۱Hilbert

^۲Integral root

^۳Touring

شدن چنین مسائلی نیست. در عمل، علاوه بر وجود الگوریتم، میزان کارآمدی^۴ الگوریتم نیز مطرح می‌گردد. به طور مثال، اگر الگوریتم حل یک مسئله مرتبه‌ی بالا یا نمایی داشته باشد، الگوریتم ارائه شده برای آن مسئله برای ورودی‌های نسبتاً بزرگ قابل اجرا نیست و نمی‌توان از آن‌ها برای حل مسئله استفاده کرد. برای تشخیص و تمیز کارآمدی الگوریتم‌های مختلف و همچنین میزان سختی مسائل در امکان ارائه‌ی الگوریتم‌های کارآمد یا غیرکارآمد، نظریه‌ی پیچیدگی^۵، دسته‌بندی‌های مختلفی برای سختی مسائل و حل‌پذیری آن‌ها ارائه داده است تا بتوان به‌طور رسمی^۶ در مورد این معیارها صحبت کرد. برای دسته‌بندی مسائل در نظریه‌ی پیچیدگی، ابتدا آن‌ها را به صورت تصمیم‌پذیر بیان می‌کنند.

مسئله‌ی ۱-۲ (مسائل تصمیم‌گیری)^۷ به دسته‌ای از مسائل گفته می‌شود که پاسخ آن‌ها تنها بله یا خیر است.

به عنوان مثال، اگر بخواهیم مسئله‌ی ۱- مرکز در فضای \mathbb{R}^d را به صورت تصمیم‌پذیر بیان کنیم، به مسئله‌ی زیر می‌رسیم:

مسئله‌ی ۲-۲ (نسخه‌ی تصمیم‌پذیر ۱- مرکز) مجموعه‌ی نقاط در فضا \mathbb{R}^d و شعاع r داده شده است، آیا دایره‌ای به شعاع r وجود دارد که تمام نقاط را بپوشاند؟

در نظریه‌ی پیچیدگی، می‌توان گفت عمده‌ترین دسته‌بندی موجود، دسته‌بندی مسائل تصمیم‌گیری به مسائل پی (P) و ان‌پی (NP) است. رده‌ی مسائل P، شامل تمامی مسائل تصمیم‌گیری است که راه‌حل چندجمله‌ای برای آن‌ها وجود دارد. از طرفی رده‌ی مسائل NP، شامل تمامی مسائل تصمیم‌گیری است که در زمان چندجمله‌ای قابل صحت‌سنجی^۸ اند. تعریف صحت‌سنجی در نظریه پیچیدگی، یعنی اگر جواب مسئله‌ی تصمیم‌گیری بله باشد، می‌توان اطلاعات اضافی با طول چندجمله‌ای ارائه داد، که در زمان چندجمله‌ای از روی آن، می‌توان جواب بله الگوریتم را تصدیق نمود. به طور مثال، برای مسئله‌ی ۱- مرکز، کافی است به عنوان تصدیق جواب بله، مرکز دایره‌ی پوشاننده، ارائه داده شود. در این صورت، می‌توان با مرتبه‌ی خطی بررسی نمود که تمام نقاط داخل این دایره قرار می‌گیرند یا نه. برای مطالعه‌ی بیش‌تر و تعاریف دقیق‌تر می‌توان به مرجع [۵] مراجعه نمود.

Efficiency^۴

Complexity theory^۵

Formal^۶

Decision problems^۷

Verifiable^۸

همان طور که می دانید درستی یا عدم درستی $P \subset NP$ از جمله معروف ترین مسائل حل نشده^۹ در نظریه پیچیدگی است. حدس بسیار قوی وجود دارد که $P \neq NP$ و بسیاری از مسائل، با این فرض حل می شوند و در صورتی که زمانی، خلاف این فرض اثبات گردد، آنگاه قسمت عمده ای از علوم کامپیوتر زیر سوال می رود.

در نظریه پیچیدگی، برای دسته بندی مسائل، یکی از روش های دسته بندی کاهش چند جمله ای^{۱۰} مسائل به یک دیگر است.

تعریف ۱-۲ می گوئیم مسئله ای A در زمان چند جمله ای به مسئله ای B کاهش می یابد، اگر وجود داشته باشد الگوریتم چند جمله ای C که به ازای هر ورودی α برای مسئله ای A ، یک ورودی β در زمان چند جمله ای برای مسئله ای B بسازد، به طوری که A ، α را می پذیرد اگر و تنها اگر B ، β را بپذیرد. در این جا منظور از پذیرفتن جواب بله به ورودی است.

از این به بعد برای سادگی به جای کاهش چند جمله ای از واژه ی کاهش استفاده می کنیم. در پی جستجو هایی که برای برابری دسته ی پی و ان پی صورت گرفت، مجموعه ای از مسائل که عمدتاً داخل ان پی هستند استخراج گردید که اگر ثابت شود یکی از آنها متعلق به پی است، آنگاه تمام مسائل دسته ی ان پی متعلق به پی خواهند بود و در نتیجه $P = NP$ خواهد بود. به این مجموعه مسائل ان پی-سخت می گویند. در واقع مسائل این دسته، مسائلی هستند که تمام مسائل داخل دسته ی ان پی، به آنها کاهش می یابند.

کوک و لوین در قضیه ای به نام **قضیه ی کوک-لوین** ثابت کردند مسئله ی صدق پذیری^{۱۱} یک مسئله ی ان پی-سخت است [۵]. با پایه قرار دادن این اثبات و استفاده از تکنیک کاهش، اثبات ان پی-سخت بودن سایر مسائل، بسیار ساده تر گردید. در ادامه مسئله ی پوشش رأسی^{۱۲} را تعریف می کنیم.

۲-۱-۱ پوشش رأسی

در این پایان نامه، از این مسئله به عنوان مسئله ی پایه برای اثبات ان پی-سخت بودن مسئله ی k -مرکز استفاده می شود. تعریف این مسئله مطابق زیر است:

^۹ Open problem

^{۱۰} Polynomial Reduction

^{۱۱} Satisfiability problem

^{۱۲} Vertex Coverage

مسئله ۲-۳ (پوشش رأسی) گراف بدون جهت $G(V, E)$ داده شده است. هدف مسئله پیدا کردن مجموعه‌ای $S \subset V$ با کمترین تعداد اعضا است به طوری که هر رأس $v \in V$ در یکی از شرایط زیر صدق کند:

$$\bullet v \in S$$

$$\bullet \text{ وجود دارد رأسی } u \in S \text{ به طوری } (v, u) \in E$$

به عبارت ساده‌تر هر رأسی یا خودش یا یکی از همسایگانش داخل مجموعه‌ای S قرار دارد.

نسخه‌ی تصمیم‌گیری این مسئله به این گونه تعریف می‌شود که آیا گراف داده‌شده دارای پوشش رأسی با اندازه‌ی k است یا نه.

قضیه ۲-۱ مسئله‌ی پوشش رأسی، یک مسئله‌ی ان‌پی-سخت است.

اثبات. برای مشاهده‌ی اثبات ان‌پی-سخت بودن مسئله‌ی پوشش رأسی، نیاز به زنجیره‌ای از مسائل که از مسئله‌ی صدق‌پذیری شروع می‌شود است، به طوری هر عضو از این زنجیره، به عضو بعدی در زمان چندجمله‌ای کاهش می‌یابد و در نهایت نتیجه می‌شود که مسئله‌ی صدق‌پذیری در زمان چندجمله‌ای به مسئله‌ی پوشش رأسی کاهش می‌یابد و در نتیجه چون مسئله‌ی صدق‌پذیری یک مسئله‌ی ان‌پی-سخت است، بنابراین مسئله‌ی پوشش رأسی نیز ان‌پی-سخت خواهد بود. برای مطالعه‌ی روند اثبات به مرجع [۵] مراجعه کنید.

□

۲-۲ الگوریتم‌های تقریبی

تا این جا با رده‌بندی مسائل به دو دسته‌ی پی و ان‌پی آشنا شدیم. نه تنها مسائل ان‌پی، بلکه بعضی از مسائل پی نیز دارای الگوریتم کارآمدی نیستند. در عمل، یک الگوریتم چندجمله‌ای با مرتبه‌ی بیش از ۳، یک الگوریتم ناکارآمد محسوب می‌شود. به طور مثال هنوز الگوریتم کارآمدی برای فهمیدن این که یک عدد اول است یا نه پیدا نشده است، با اینکه الگوریتم ارائه شده یک الگوریتم چندجمله‌ای است.

عمده‌ی مسائل کاربردی مطرح در دنیای واقع، یا متعلق به دسته‌ی ان‌پی هستند و در نتیجه راه‌حل چندجمله‌ای ندارند، یا اگر راه‌حل چند جمله‌ای داشته باشند، مرتبه‌ی چندجمله‌ای بالاست و در نتیجه راه‌حل کارآمدی محسوب نمی‌گردد. یکی از رویکردهای رایج در برابر چنین مسائلی، صرف نظر کردن از دقت راه‌حل‌هاست. به طور مثال راه‌حل‌های مکاشفه‌ای^{۱۳} گوناگونی برای مسائل مختلف به خصوص مسائل ان‌پی بیان شده است. این راه‌حل‌ها بدون این‌که تضمین کنند راه‌حل خوبی ارائه می‌دهند یا حتی جوابشان به جواب بهینه نزدیک است، اما با معیارهایی سعی در بهینه عمل کردن دارند و تا حد ممکن سعی در ارائه‌ی جواب بهینه یا نزدیک بهینه دارند. اما در عمل، تنها برای دسته‌ای از کاربردها پاسخ قابل قبولی ارائه می‌دهند.

مشکل عمده‌ی راه‌حل‌های مکاشفه‌ای، عدم امکان استفاده از آن‌ها برای تمام کاربردها است. بنابراین در رویکرد دوم که اخیراً نیز مطرح شده است و عمر کمی دارد، سعی در ارائه‌ی الگوریتم‌های مکاشفه‌ای شده است که تضمین می‌کنند اختلاف زیادی با الگوریتمی که جواب بهینه می‌دهد، نداشته باشند. در واقع این الگوریتم‌ها همواره و در هر شرایطی، تقریبی از جواب بهینه را ارائه می‌دهند. به چنین الگوریتم‌هایی، **الگوریتم‌های تقریبی**^{۱۴} می‌گویند. علت اصلی این نام‌گذاری، تقریب زدن جواب الگوریتم بهینه است. ضریب تقریب یک الگوریتم تقریبی، به حداکثر نسبت جواب الگوریتم تقریبی به جواب بهینه گفته می‌شود.

الگوریتم‌های تقریبی تنها به علت محدودیت کارایی الگوریتم‌هایی که جواب بهینه می‌دهند، مورد استفاده قرار نمی‌گیرند. هر نوع محدودیتی ممکن است، استفاده از الگوریتم تقریبی را نسبت به الگوریتمی که جواب بهینه می‌دهد، مقرون به صرفه کند. به طور مثال از جمله عوامل دیگری که ممکن است باعث این انتخاب شود، کاهش میزان حافظه‌ی مصرفی باشد. برای طیف وسیعی از مسائل، کمبود حافظه، باعث می‌شود الگوریتم‌هایی با حافظه‌ی مصرفی کمتر طراحی شود که به دقت الگوریتم‌های بهینه عمل نمی‌کند اما می‌تواند با مصرف کمتر به تقریبی از جواب بهینه دست یابد. معمولاً چنین الگوریتم‌هایی حافظه‌ی مصرفی از مرتبه‌ی زیرخطی^{۱۵} دارند و به همین دلیل برای داده‌های حجیم بسیار کاربرد دارند.

^{۱۳} Heuristic^{۱۴} Approximation Algorithm^{۱۵} Sublinear

مسئله	کران پایین تقریب پذیری
پوشش رأسی	$1,3606 [6]$
k -مرکز	$2 [7]$
k -مرکز در فضای اقلیدسی	$1,822 [8]$
۱-مرکز در حالت جویبار داده	$\frac{1+\sqrt{2}}{4} [9]$
k -مرکز با نقاط پرت و نقاط اجباری	$3 [3]$

جدول ۲-۱: نمونه‌هایی از ضرایب تقریب برای مسائل بهینه‌سازی

۲-۲-۱ میزان تقریب پذیری مسائل

همان‌طور که تا این‌جا دیدیم، یکی از راه‌کارهایی که برای کارآمد کردن راه‌حل ارائه شده برای یک مسئله وجود دارد، استفاده از الگوریتم‌های تقریبی برای حل آن مسئله است. یکی از عمده‌ترین دغدغه‌های مطرح در الگوریتم‌های تقریبی کاهش ضریب تقریب است. حتی در بعضی از موارد حتی امکان ارائه‌ی الگوریتم تقریبی با ضریبی ثابت نیز وجود ندارد. به‌طور مثال، همان‌طور که در فصل کارهای پیشین بیان خواهد شد، الگوریتم تقریبی با ضریب تقریب کم‌تر از ۲، برای مسئله‌ی k -مرکز وجود ندارد مگر اینکه $P = NP$ باشد. برای مسائل مختلف، معمولاً می‌توان کران پایینی برای میزان تقریب پذیری آن‌ها ارائه داد. در واقع برای مسائل ان‌پی، علاوه بر این الگوریتم کارآمدی وجود ندارد، بعضاً الگوریتم تقریبی برای حل آن با ضریبی تقریب کم و نزدیک به یک نیز وجود ندارد. در جدول ۲-۱ از میزان تقریبی مسائل مختلفی که در این پایان‌نامه مورد استفاده قرار می‌گیرد ببینید.

۲-۳ الگوریتم‌های جویبار داده

در علوم کامپیوتر، الگوریتم جویبار داده به الگوریتم‌هایی گفته می‌شود که برای پردازش جویبارهای داده طراحی شده‌اند به طوری که ورودی آن، به صورت دنباله‌ای از داده‌ها داده می‌شود و تنها می‌توان تعدادی محدود بار، از روی دنباله گذر کرد (معمولاً تنها یک بار). الگوریتم‌های جویبار داده معمولاً محدودیت شدیدی در حافظه دارند (نسبت به اندازه‌ی ورودی) و به علت تعداد زیاد داده‌ها، محدودیت زمانی پردازش برای هر داده نیز مطرح است.

چنین محدودیت‌هایی معمولاً باعث می‌شود که الگوریتم جویبار داده تنها بتواند یک جواب تقریبی از جواب بهینه را با استفاده از اطلاعات مختصری که در حافظه نگه می‌دارد را ارائه دهد.

۲-۳-۱ مقدمه

در سال‌های اخیر، پیشرفت‌های حوزه‌ی تکنولوژی، امکان جمع‌آوری داده‌ها را به صورت پوسته ممکن ساخته است. به طور مثال می‌توان از تراکنش‌های بانکی، استفاده از تلفن همراه یا مرورگر وب خود را در نظر بگیرید. در هر تعاملی که با این سیستم‌ها انجام می‌دهید، میزان زیادی داده تولید شده و ذخیره می‌گردد. در اکثر مواقع می‌توان این حجم عظیم داده را مورد پردازش قرار داد و اطلاعات بسیار مفیدی از آن، استخراج نمود. زمانی که حجم داده بسیار زیاد باشد، معمولاً چالش‌های محاسباتی و الگوریتمی برای الگوریتم‌ها به وجود می‌آورد. از جمله‌ی آن‌ها می‌توان به موارد زیر اشاره کرد:

- با افزایش حجم داده، امکان پردازش داده‌ها به صورت کارآمد با چندبار عبور کردن از جویبار داده وجود ندارد. یکی از مهم‌ترین موانع در طراحی الگوریتم‌های کارآمد برای مدل جویبار داده این است که الگوریتم‌ها مجبور هستند هر داده را حداکثر یک بار مورد پردازش قرار دهند. بنابراین الگوریتم‌های مدل جویبار داده، معمولاً تک‌گذره‌اند.
- در اکثر الگوریتم‌های جویبار داده، از یک واحد برای پردازش داده‌ها به صورت محلی استفاده می‌شود. علت اصلی وجود چنین واحدی، نیاز این الگوریتم‌ها به تشخیص تغییرات در داده‌های ورودی است. این عملکرد الگوریتم‌های جویبار داده را می‌توان نوعی استفاده از اصل محل‌گرایی^{۱۶} محسوب نمود. اما مشکل عمده‌ی این نوع رویکرد اجتناب ناپذیر، در عمده‌ی موارد، عدم امکان راه‌حلی مناسب برای مسئله است. در ساخت الگوریتم‌های جویبار داده، باید دقت و تمرکز عمده‌ای را صرف تشخیص نحوه‌ی تغییر داده‌های ورودی نمود.
- یکی از مهم‌ترین مشخصه‌های جویبار داده‌ها، ماهیت عدم متمرکز بودن داده‌ها است. از طرفی یک پردازنده به تنهایی دارای محدودیت بسیار زیادی از نظر قدرت پردازشی و حافظه است. بنابراین معمولاً الگوریتم‌های جویبار داده، به گونه‌ای طراحی می‌شوند که قابل توزیع‌پذیری بوده و توانایی اجرا شدن به صورت چندروندی^{۱۷} را دارا باشند.

^{۱۶} Locality

^{۱۷} Multi Thread

۲-۳-۲ گونه‌های مطرح

در مدل جویبار داده، بعضی یا همه‌ی نقاط ورودی که باید پردازش گردند برای دسترسی تصادفی^{۱۸} از حافظه‌ی اصلی یا حافظه‌ی جانبی در دسترس نیستند، بلکه به مرور زمان، به صورت جویبار یا جویبارهایی از داده در اختیار الگوریتم قرار می‌گیرند [۲].

هر جویبار را می‌توان به عنوان دنباله‌ای مرتب از نقاط (یا به‌روزرسانی‌ها^{۱۹}) در نظر گرفت به طوری که به ترتیب دنباله قابل دسترسی هستند و هر کدام را تنها به تعدادی محدود بار (معمولاً یک بار) می‌توان خواند [۲].

مسائل زیادی در این مدل مورد بررسی قرار گرفته‌اند که از جمله مهم‌ترین آن‌ها می‌توان به موارد زیر اشاره کرد:

- محاسبه‌ی آماری مربوط به توزیع داده‌های یک جویبار داده که به قدری بزرگ هستند قابل ذخیره شدن نیست.
- مسائل مربوط به گراف‌ها، به طوری که ماتریس مجاورت گراف به ترتیب دلخواهی به صورت جویبار داده به الگوریتم داده می‌شود.
- مسائل مربوط به دنباله‌ها و استخراج اطلاعات از دنباله‌ها که وابستگی شدیدی به ترتیب اعضای دنباله دارند، مانند پیدا کردن طولانی‌ترین زیردنباله با اعضای صعودی یا پیدا کردن تعداد نابه‌جایی‌های^{۲۰} داخل دنباله. [۲]

۲-۳-۳ تحلیل الگوریتم‌های جویبار داده

کارایی یک الگوریتم جویبار داده بر اساس سه معیار زیر اندازه‌گیری می‌شود:

- تعداد مرتبه‌ای که الگوریتم از روی جویبار داده گذر می‌کند
- میزان حافظه‌ی مصرفی
- زمان مصرفی به ازای هر داده

^{۱۸} Random access

^{۱۹} Update

^{۲۰} Inversion

الگوریتم‌های جویبار داده تشابه‌های زیادی با الگوریتم‌های برخط^{۲۱} دارند. به طور مثال، هر دو نوع الگوریتم، نیاز دارند قبل از اینکه تمام ورودی فرا برسد، در مورد داده‌ی جویبار داده تصمیم بگیرند. اما تفاوت‌های عمده‌ای نیز در این الگوریتم‌ها وجود دارد. الگوریتم‌های جویبار داده، دارای حافظه‌ی محدودی هستند، اما می‌توانند تصمیم‌گیری راجع به یک داده را تا چند مرحله به عقب بیاندازند، در صورتی که الگوریتم‌های برخط، به محض ورود داده، باید در مورد آن تصمیم بگیرند.

۲-۳-۴ مجموعه هسته

یکی از تکنیک‌های رایج در الگوریتم‌های جویبار داده، نگه‌داری نماینده‌ای با اندازه‌ی بسیار کوچک‌تر نسبت جویبار داده است. این مجموعه معمولاً دغدغه‌ی محدودیت حافظه را در الگوریتم‌های جویبار داده برطرف می‌کند. به چنین مجموعه‌ای، مجموعه هسته می‌گوییم. حال به تعریف رسمی هسته می‌پردازیم:

تعریف ۲-۲ فرض کنید μ یک تابع اندازه‌گیری^{۲۲} (همانند تابع عرض مجموعه‌ای از نقاط) از \mathbb{R}^d به اعداد حقیقی نامنفی $\mathbb{R}^d \cup \{0\}$ باشد. فرض کنید که این تابع، یک تابع یکنواخت است، یعنی به ازای هر دو مجموعه‌ی P_1 و P_2 که $P_1 \subset P_2$ است، آنگاه

$$\mu(P_1) \leq \mu(P_2)$$

فرض کنید $\epsilon > 0$ داده شده است، به زیرمجموعه‌ی $Q \subset P$ یک ϵ -مجموعه‌ی هسته برای مجموعه P است اگر رابطه‌ی زیر برقرار باشد:

$$(1 - \epsilon)\mu(P) \leq \mu(Q)$$

یکی از مجموعه هسته‌های معروف، مجموعه هسته‌ی مطرح برای برای تابع اندازه‌گیری عرض نقاط است. به چنین مجموعه هسته‌ای به اختصار ϵ -هسته^{۲۳} گفته می‌شود. تعریف دقیق ϵ -هسته در زیر آمده است:

تعریف ۲-۳ فرض کنید S^{d-1} کره‌ی واحد با مرکز واقع در مبدأ در فضای \mathbb{R}^d باشد. به ازای هر مجموعه P از نقاط در فضای \mathbb{R}^d و هر جهت دلخواه $u \in S^{d-1}$ ، عرض جهت‌دار P در جهت u که با نماد $\omega(u, P)$

^{۲۱}Online

^{۲۲}Measure Function

^{۲۳} ϵ -Kernel

تعریف می‌شود، مطابق زیر است:

$$\omega(u, P) = \max_{p \in P} \langle u, p \rangle - \min_{p \in P} \langle u, p \rangle$$

که در آن $\langle \cdot, \cdot \rangle$ همان ضرب داخلی دو بردار است. برای متغیر $\epsilon > 0$ ، زیرمجموعه‌ی $Q \subset P$ یک ϵ -هسته نامیده می‌شود اگر به ازای هر $u \in S^{d-1}$ داشته باشیم:

$$(1 - \epsilon)\omega(u, P) \leq \omega(u, Q)$$

ϵ -هسته یکی از اساسی‌ترین مجموعه‌های هسته‌های مطرح است و برای طیف وسیعی از مسائل قابل استفاده است. الگوریتم‌های زیادی برای محاسبه‌ی ϵ -هسته در حالت ایستا تعریف شده است [۱۰].

یکی از روش‌هایی که در تبدیل یک الگوریتم از حالت ایستا به حالت جویبار داده استفاده از روش دوبرابری‌سازی^{۲۴} است. این روش یک روش عام و قابل استفاده برای طیف وسیعی از مسائل است. به عنوان نمونه، دوبرابری‌سازی را بر روی ϵ -هسته مورد بررسی قرار می‌دهیم. ϵ -هسته، دارای دو خاصیت اساسی است که آن را قابل استفاده برای تکنیک دوبرابری‌سازی است.

• اگر P_2 یک ϵ -هسته برای P_1 باشد و P_3 یک δ -هسته برای P_2 باشد، آنگاه P_3 یک $(\epsilon + \delta)$ -هسته برای P_1 است.

• اگر P_2 یک ϵ -هسته برای P_1 باشد و Q_2 یک ϵ -هسته برای Q_1 باشد، آنگاه $P_2 \cup Q_2$ یک ϵ -هسته برای $P_1 \cup Q_1$ است.

ایده‌ی اصلی تکنیک دوبرابری‌سازی، استفاده از روش مبنای دو است، که در ابتدا از نقطه‌ی اول یک مجموعه هسته ساخته می‌شود. سپس با نقطه‌ی بعدی یک مجموعه هسته از دو نقطه‌ی فعلی ساخته می‌شود. سپس با دو نقطه‌ی بعدی یک مجموعه‌ی هسته از چهار نقطه فعلی ساخته می‌شود. اگر این روند را به صورت توان‌هایی از ۲ ادامه بدهیم، در هر لحظه حداکثر، $\log(n)$ مجموعه هسته داریم که در طول الگوریتم، بعضی از آن‌ها با یک‌دیگر ادغام می‌گردند و پس از ادغام، یک هسته از کل آن‌ها به عنوان نماینده انتخاب می‌گردد و مجموعه هسته‌ی جدیدی از روی آن ساخته می‌شود. در واقع اگر دوباره بر روی دو مجموعه‌ی ادغام شده هسته حساب نشود، اندازه‌ی هسته به صورت نمایی افزایش می‌یابد که مطلوب نیست. ازین رو برای کاهش حافظه‌ی مصرفی از روی دو هسته‌ی ادغام شده، یک هسته‌ی

^{۲۴}Doubling

جدید محاسبه می‌گردد. این کار، ضریب تقریب را افزایش می‌دهد ولی باعث می‌شود، حافظه‌ی مصرفی، حداکثر $\log(n)$ برابر حافظه‌ی مصرفی در حالت ایسا می‌شود.

فصل ۳

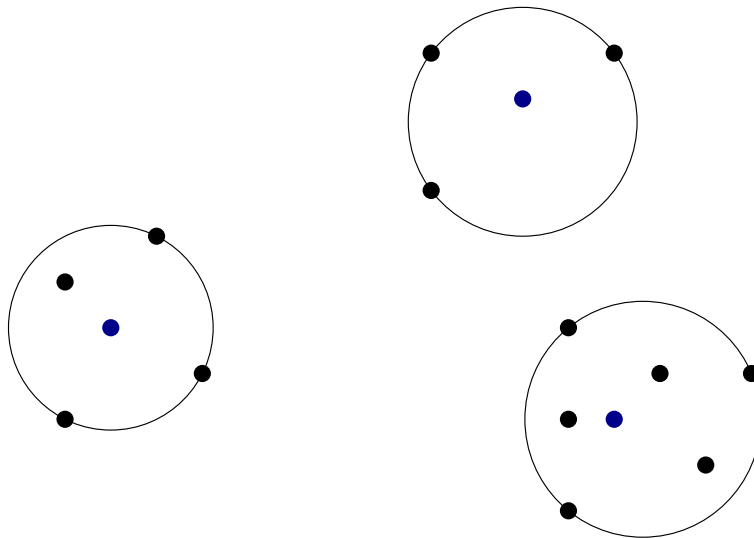
کارهای پیشین

در این فصل کارهای پیشین انجام شده روی مسئله k - مرکز، در سه بخش مورد بررسی قرار می‌گیرد. در بخش اول، مسئله k - مرکز مورد بررسی قرار می‌گیرد. در بخش دوم، حالت جویبار داده‌ی مسئله و مجموعه هسته‌های مطرح برای این مسئله مورد بررسی قرار می‌گیرد. در نهایت، در بخش سوم، مسئله k - مرکز با داده‌های پرت مورد بررسی قرار می‌گیرد.

۳-۱ k - مرکز در حالت ایستا

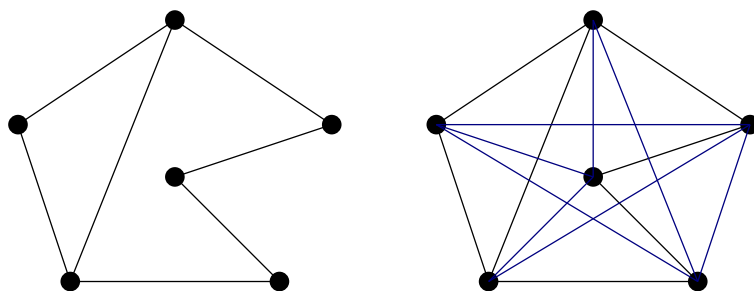
مسئله k - مرکز به عنوان مسئله‌ی شناخته شده در علوم کامپیوتر مطرح است. این مسئله، در واقع یک مسئله‌ی بهینه‌سازی است که سعی در کاهش بیش‌ترین فاصله نقاط از مرکز دسته‌ها را دارد. سختی اصلی این مسئله در انتخاب مرکز دسته‌هاست. زیرا اگر بتوانیم مرکز دسته‌ها را به درستی تشخیص دهیم، کافی است هر نقطه را به دسته‌ای که نزدیک‌ترین مرکز را دارد، تخصیص دهیم. به وضوح چنین تخصیصی بهینه‌ترین تخصیص ممکن است. نمونه‌ای از این تخصیص را در شکل ۳-۱ نشان داده شده است.

در سال ۱۹۷۹، نه تنها اثبات گردید که این مسئله در حالت کلی یک مسئله‌ی ان‌پی-سخت است [۱۱]، بلکه ثابت شده است که این مسئله در صفحه‌ی دو بعدی و با متریک اقلیدسی نیز ان‌پی-سخت است [۱۲]. فراتر از این، ثابت شده است که برای مسئله k - مرکز با متریک دلخواه هیچ الگوریتم تقریبی با ضریب تقریب بهتر از ۲ وجود ندارد. ایده‌ی اصلی این کران پایین، کاهش مسئله‌ی پوشش رأسی، به مسئله k - مرکز است. برای چنین کاهش‌ی کافی است، از روی گراف اصلی که می‌خواهیم



شکل ۳-۱: نمونه‌ای از تخصیص نقاط به ازای مراکز آبی رنگ.

کوچک‌ترین مجموعه‌ی پوشش رأسی را در آن پیدا کنیم، یک گراف کامل بسازیم به طوری که به ازای هر یال در گراف اصلی، یک یال با وزن یک و به ازای هر یال که در گراف اصلی وجود ندارد، یک یال با وزن ۲ قرار می‌دهیم. نمونه‌ای از چنین تبدیلی را در شکل ۳-۲ می‌توانید مشاهده کنید. حال اگر الگوریتمی بتواند مسئله‌ی k -مرکز را با ضریب تقریب بهتر از ۲ حل نماید، آنگاه گراف جدید دارای یک k -مرکز با شعاع کمتر از ۲ است، اگر و تنها اگر گراف اصلی دارای یک پوشش رأسی با اندازه‌ی k باشد. برای متریک L_2 یا فضای اقلیدسی^۱ نیز ثابت شده است برای مسئله‌ی k -مرکز، الگوریتم تقریبی با ضریب تقریب بهتر از $1/822$ وجود ندارد [۸].



شکل ۳-۲: نمونه‌ای از تبدیل یک گراف ورودی مسئله‌ی پوشش رأسی به یک ورودی مسئله‌ی k -مرکز (در گراف سمت چپ، یال‌های سیاه وزن ۱ و یال‌های آبی، وزن ۲ دارند)

^۱Euclidean space

الگوریتم ۱ گزنالز

- ۱: S را برابر مجموعه تهی قرار بده.
- ۲: عنصر دلخواه از مجموعه نقاط V را به S اضافه کن.
- ۳: به ازای i بین ۲ تا k :
- ۴: v را نقطه‌ای از V در نظر بگیرید که بیشترین فاصله را از مجموعه‌ی S دارد.
- ۵: v را S اضافه کن.
- ۶: S را برگردان

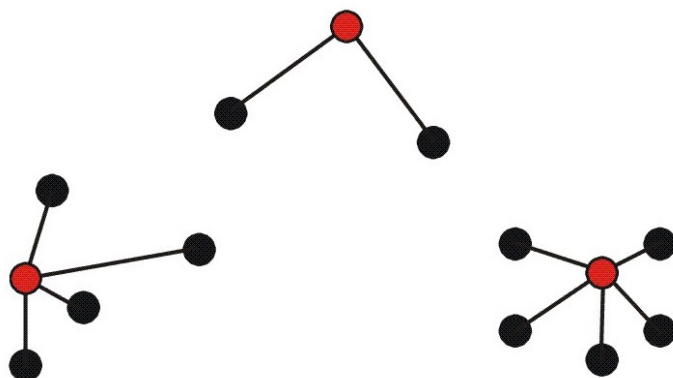
یکی از اولین الگوریتم‌های تقریبی برای مسئله‌ی k -مرکز به وسیله‌ی گزنالز^۲ ارائه شده است [۱۳]. این الگوریتم یک الگوریتم تقریبی با ضریب ۲ است و در زمان $O(kn)$ قابل اجراست. الگوریتم گزنالز، از نظر روش برخورد با مسئله، یک الگوریتم حریصانه^۳ محسوب می‌شود. برای این‌که عملکرد الگوریتم گزنالز را درک کنید، نیاز به تعریف فاصله‌ی یک نقطه از یک مجموعه نقطه داریم.

تعریف ۳-۱ فاصله‌ی نقطه‌ی v از مجموعه‌ای ناتهی از نقاط S را برابر فاصله‌ی نقطه‌ای درون S از v تعریف می‌کنیم، به گونه‌ای که از تمام نقاط S به v نزدیک‌تر باشد. در واقع داریم:

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

همان‌طور که در الگوریتم ۱ مشاهده می‌کنید، روش اجرای این الگوریتم به این گونه است که در ابتدا یک نقطه‌ی دلخواه را به عنوان مرکز دسته‌ی اول در نظر می‌گیرد. سپس دورترین نقطه از آن را به عنوان مرکز دسته‌ی دوم در نظر می‌گیرد. در هر مرحله، دورترین نقطه از مرکز مجموعه دسته‌های موجود را به عنوان مرکز دسته‌ی جدید به مجموعه مراکز دسته‌ها اضافه می‌کند. با اجرای الگوریتم تا k مرحله، مراکز دسته‌ها انتخاب می‌شود. حال اگر هر نقطه را به نزدیک‌ترین مرکز انتخابی تخصیص دهیم، می‌توان نشان داد که شعاع بزرگ‌ترین دسته، حداکثر دو برابر شعاع بهینه برای مسئله‌ی k -مرکز است. فدر^۴ و سایرین، زمان اجرای الگوریتم گزنالز را برای هر L_p -متریک به مرتبه‌ی $O(n \log k)$ بهبود بخشیدند. نمونه‌ای از اجرای الگوریتم گزنالز، در شکل ۳-۳ نشان داده شده است.

^۲Gonzalez^۳Greedy^۴Feder



شکل ۳-۳: نمونه‌ای از حل مسئله‌ی ۳-مرکز با الگوریتم گنزالز

تا به اینجا، تنها بر روی حالت کلی مسئله‌ی k -مرکز صحبت شد و تنها محدودیتی که مورد توجه قرار گرفت، متریک مطرح برای فاصله‌ی نقاط بوده است. در ادامه به بررسی، حالاتی از مسئله‌ی k -مرکز، که k تعداد دسته‌ها یا d ابعاد فضا ثابت باشند می‌پردازیم. آگاروال^۵ و سایرین الگوریتمی دقیق با زمان اجرای $n^{O(k^{1-\frac{1}{d}})}$ برای مسئله k -مرکز در فضای L_p -متریک با ابعاد ثابت d ارائه داده‌اند [۱۴]. قابل توجه است که اگر d ثابت نباشد، مسئله‌ی k -مرکز حتی برای متریک اقلیدسی (L_2 -متریک) با تعداد دسته‌ی ثابت $k \geq 2$ ، ان‌پی-سخت است [۱۵].

علاوه بر حالاتی که ابعاد فضا یا تعداد دسته‌ها ثابت‌اند، مسئله‌ی k -مرکز برای حالتی که مقادیر k و d کوچک هستند، مورد بررسی قرار گرفته‌اند و الگوریتم‌های بهینه‌تری از الگوریتم‌های کلی برای این حالت‌های خاص ارائه شده است. به طور مثال، برای مسئله‌ی ۱-مرکز در فضای اقلیدسی با ابعاد ثابت، الگوریتم خطی با زمان اجرای $O((d+1)n)$ وجود دارد [۱۶]. الگوریتم ارائه شده بر پایه‌ی دو نکته‌ی اساسی بنا شده است. اول اینکه کره‌ی بهینه را می‌توان با حداکثر $d+1$ نقطه‌ی واقع در پوسته‌ی کره‌ی بهینه مشخص نمود و دوم اینکه اگر نقاط ورودی را با ترتیبی تصادفی پیمایش کنیم احتمال اینکه نقطه‌ی پیمایش شده جزء نقاط مرزی باشد از مرتبه‌ی $O(\frac{d}{n})$ است که با توجه به ثابت بودن d این احتمال برای n های بزرگ کوچک محسوب می‌شود و زمان اجرای الگوریتم، به طور میانگین خطی خواهد بود.

در صفحه‌ی اقلیدسی برای مسئله‌ی ۲-مرکز، بهترین الگوریتم را چن^۶ با زمان اجرای $O(c \log^2 n \log^2 \log n)$ و حافظه‌ی $O(n)$ ارائه داده است [۱۷]. برای فضای سه‌بعدی اقلیدسی نیز آگاروال و سایرین، الگوریتمی

Agarwal^۵
Chan^۶

با متوسط زمان اجرای $O(n^3 \log^k n)$ ارائه داده است [۱۸].

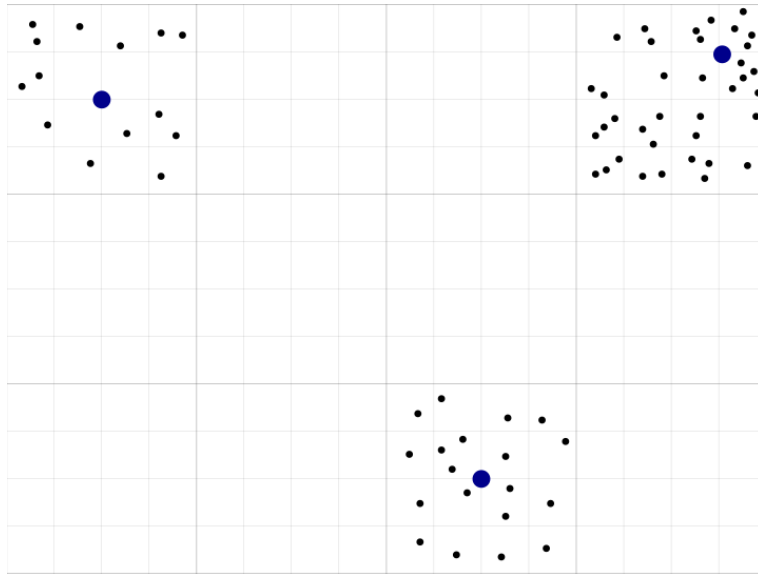
۳-۲ k -مرکز در حالت جویبار داده

در مدل جویبار داده، مشکل اصلی عدم امکان نگه‌داری تمام داده‌ها در حافظه است و باید سعی شود تنها داده‌هایی که ممکن است در ادامه مورد نیاز باشد را، نگه‌داریم. یکی از راه‌های رایج برای این کار نگه‌داری مجموعه‌ای از نقاط (نه لزوماً زیرمجموعه‌ای از نقاط ورودی) به عنوان نماینده‌ی نقاط به طوری که جواب مسئله‌ی k -مرکز برای آن‌ها منطبق با جواب مسئله‌ی k -مرکز برای نقاط اصلی باشد (با تقریب قابل قبولی). به چنین مجموعه‌ای مجموعه‌ی هسته‌ی نقاط گفته می‌شود.

بهترین مجموعه هسته‌ای که برای مسئله‌ی k -مرکز ارائه شده است، روش ارائه‌شده به وسیله‌ی ضربایی‌زاده برای نگه‌داری یک ϵ -هسته با حافظه‌ی $O(\frac{k}{\epsilon^d})$ برای L_p -متریک‌ها است [۱۹]. در روش ارائه شده، از چند ایده‌ی ترکیبی استفاده شده است. در ابتدا، الگوریتم با استفاده از یک الگوریتم تقریبی با ضریب ثابت، یک تقریب از جواب بهینه به دست می‌آورد. به طور مثال با استفاده از الگوریتم گنزالز، یک دو تقریب از شعاع بهینه به علاوه‌ی مرکز دسته‌های پیدا شده را به دست می‌آورد. حال کافی است که با طول شعاع الگوریتم تقریبی، حول هر مرکز به دست آمده، یک توری با $O(\frac{1}{\epsilon})$ شبکه‌بندی در هر بعد تشکیل داد و چون هر نقطه در حداقل یکی از توری‌ها قرار می‌گیرد، می‌توان با حداکثر ϵ تقریب در جواب نهایی نقاط را به نقاط شبکه‌بندی توری گرد نمود. با این کار، دیگر ما نیازی به نگهداری تمام نقاط ورودی نداشته و تنها نیاز به نگهداری نقاط شبکه‌بندی توری داریم. با این روش می‌توان به یک ϵ -هسته برای مسئله‌ی k -مرکز رسید. نکته‌ی اساسی برای سازگار سازی روش ارائه‌شده با مدل جویبار داده‌ی تک‌گذره استفاده از روش دوبرارسازی^۷ رایج در الگوریتم‌های جویبار داده است. نمونه‌ای از اجرای الگوریتم ضربایی‌زاده را در شکل ۳-۴ نشان داده شده است. برای دیدن اثبات‌ها و توضیح بیش‌تر در مورد روش ارائه شده می‌توانید به مرجع [۱۹] مراجعه کنید.

از جمله مشکلات وارده به الگوریتم ضربایی‌زاده، وابستگی اندازه‌ی مجموعه‌ی هسته به ابعاد فضا است. بنابراین هسته‌ی ارائه شده به وسیله‌ی ضربایی‌زاده را نمی‌توان برای ابعاد بالا مورد استفاده قرار داد. از طرفی، حساب کردن جواب از روی هسته در زمان چندجمله‌ای بر اساس k و d قابل انجام نیست. با توجه با موارد گفته شده، برای قابل استفاده شدن الگوریتم‌ها برای ابعاد بالا، الگوریتم‌هایی

^۷Doubling



شکل ۳-۴: نمونه‌ای از توری‌بندی الگوریتم ضربی‌زاده (نقاط آبی، مراکز به دست آمده از الگوریتم تقریبی است). پس از توری‌بندی کافی است برای هر کدام از خانه‌های شبکه‌بندی، تنها یکی را به عنوان نماینده در نظر بگیریم.

ارائه می‌شود که ضریب تقریب بدتری دارند، اما میزان حافظه‌ی مصرفی یا اندازه‌ی مجموعه هسته‌ی آن‌ها چندجمله‌ای بر اساس d و k و $\log n$ باشد. به چنین الگوریتم‌هایی الگوریتم‌های جویبار داده برای ابعاد بالا گفته می‌شود.

اولین الگوریتم ارائه شده برای ابعاد بالا، الگوریتمی با ضریب تقریب ۸ و حافظه‌ی مصرفی $\mathcal{O}(dk)$ است [۲۰]. پس از آن، گوها^۸، به طور موازی با مک‌کاتن^۹ و سایرین، الگوریتمی با ضریب تقریب $(2 + \epsilon)$ با حافظه‌ی $\mathcal{O}(\frac{dk}{\epsilon} \log \frac{1}{\epsilon})$ برای مسئله‌ی k -مرکز در هر فضای متریکی ارائه دادند [۴، ۲۱]. در سال ۲۰۱۴، اهن^{۱۰} و سایرین، الگوریتمی با همین ضریب تقریب و حافظه‌ی $\mathcal{O}((k+3)! 2^{k \frac{d}{\epsilon}})$ ارائه داده‌اند که برای k های ثابت، حافظه را از مرتبه‌ی $\mathcal{O}(\log \frac{1}{\epsilon})$ کاهش می‌دهد [۲۲].

تا به این‌جا ما به بررسی مسئله‌ی k -مرکز در حالت جویبار داده بدون محدودیت خاصی پرداختیم. برای حالت‌های خاص k ، به خصوص ۱ و ۲، مسئله‌ی k -مرکز با متریک اقلیدسی مورد توجه زیادی قرار گرفته است و راه‌حل‌های بهینه‌تری نسبت به حالت کلی برای آن‌ها پیشنهاد شده است. به طور مثال،

Guha^۸McCutchen^۹Ahn^{۱۰}

می‌توان یک هسته با اندازه‌ی $O(\frac{1}{\epsilon^{\frac{d-1}{2}}})$ ، با استفاده از نقاط حدی^{۱۱} در تعداد مناسبی جهت، به صورت جویبار داده ساخت.

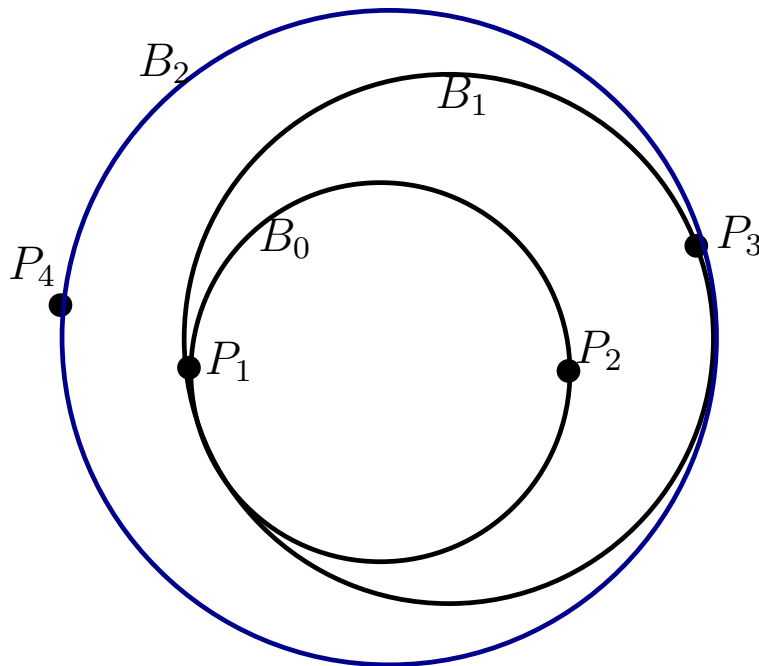
الگوریتم ۲ الگوریتم ضربی‌زاده

- ۱: B را تویی به مرکز نقطه‌ی اول و شعاع صفر قرار دهید.
- ۲: به ازای هر نقطه‌ی u در جویبار داده:
- ۳: اگر u داخل B قرار می‌گیرد:
- ۴: ادامه بده.
- ۵: در غیر این صورت:
- ۶: B را با کوچک‌ترین توپ ممکن که هردوی B و u را می‌پوشاند، جایگزین کن.
- ۷: مجموعه‌ی S را برگردان.

همان‌طور که برای الگوریتم ضربی‌زاده مطرح شد، مشکل عمده‌ی مجموعه هسته‌ی ارائه‌شده، وابستگی حافظه‌ی مصرفی آن به d است. در راستای حل این مشکل، ضربی‌زاده و سایرین [۲۳] برای ابعاد بالا و متریک اقلیدسی، الگوریتمی با ضریب تقریب $1/5$ و حافظه‌ی مصرفی $O(d)$ ارائه دادند. در واقع در الگوریتم آن‌ها، در هر لحظه تنها یک مرکز و یک شعاع را نگه داشته می‌شود، که کم‌ترین حافظه‌ی ممکن برای مسئله‌ی ۱-مرکز است. همان‌طور که در الگوریتم ۲ مشاهده می‌کنید، نقطه‌ی اول را به عنوان مرکز کره با شعاع صفر در نظر گرفته می‌شود. با فرا رسیدن هر نقطه‌ی جدید، اگر نقطه‌ی مد نظر، داخل کره‌ی فعلی بیفتد که بدون هیچ تغییری ادامه داده و در صورتی که بیرون کره‌ی فعلی بیفتد، آن را با کوچک‌ترین کره‌ای که نقطه‌ی جدید به علاوه‌ی کره‌ی قبلی را به طور کامل می‌پوشاند، جایگزین می‌کند.

به وضوح در هر لحظه کره‌ی ساخته شده تمام نقاطی که تا کنون در جویبار داده آمده است را می‌پوشاند. از طرفی ثابت می‌شود شعاع کره در هر لحظه، حداکثر $1/5$ - برابر شعاع کره‌ی بهینه است. نکته‌ی قابل توجه در مورد این الگوریتم، نگهداری کم‌ترین حافظه‌ی ممکن برای مسئله‌ی ۱-مرکز است. زیرا تنها یک مرکز و یک شعاع نگه می‌دارد. نمونه‌ای از اجرای الگوریتم را بر روی چهار نقطه می‌توان در شکل ۳-۵ دید. برای اثبات کامل‌تر می‌توانید به مرجع [۲۳] مراجعه کنید.

^{۱۱}Extreme points



شکل ۳-۵: نمونه‌ای از اجرای الگوریتم ضربی‌زاده بر روی چهار نقطه $P_1 \dots P_4$ که به ترتیب اندیس در جویبار داده فرا می‌رسند و دایره‌های B_0, \dots, B_2 دایره‌هایی که الگوریتم به ترتیب نگه می‌دارد.

در ادامه، آگاروال و سایرین [۹] الگوریتمی تقریبی با حافظه‌ی مصرفی $O(d)$ ارائه دادند. در الگوریتم ارائه شده، ضریب تقریب، برابر $\frac{1+\sqrt{3}}{4}$ تخمین زده شد، اما با تحلیل دقیق‌تری که چن و سایرین [۲۴] بر روی همان الگوریتم انجام دادند، مشخص شد همان الگوریتم دارای ضریب تقریب $1/22$ است.

الگوریتم آگاروال از الگوریتم کلارکسون و سایرین [۲۵] به عنوان یک زیرالگوریتم استفاده می‌کند. الگوریتم کلارکسون، الگوریتمی کاملاً مشابه الگوریتم گنزالز است و به این گونه عمل می‌کند که در ابتدا یک نقطه دلخواه را به عنوان نقطه‌ی اول انتخاب می‌کند، سپس دورترین نقطه از نقطه‌ی اول را به عنوان دومین نقطه انتخاب می‌کند. ازین به بعد در هر مرحله، نقطه‌ای که از نقاط انتخاب شده‌ی قبلی بیش‌ترین فاصله را دارد به عنوان نقطه‌ی جدید انتخاب می‌کند. اگر این الگوریتم را تا $O(\frac{1}{\epsilon})$ مرحله ادامه بدهیم، به مجموعه‌ای با اندازه‌ی $O(\frac{1}{\epsilon})$ خواهیم رسید که کلارکسون و سایرین اثبات کرده‌اند که یک ϵ -هسته برای مسئله‌ی ۱-مرکز است. در واقع آن‌ها نشان داده‌اند یک $O(\frac{1}{\epsilon})$ -مرکز به دست آمده برای مجموعه‌ای از نقاط، یک ϵ -هسته برای مسئله‌ی ۱-مرکز برای همان مجموعه نقاط است.

الگوریتم آگاروال به این گونه عمل می‌کند که اولین نقطه‌ی جویبار داده را به عنوان تنها نقطه‌ی مجموعه‌ی K_1 در نظر می‌گیرد. حال تا وقتی که نقاطی که فرا می‌رسند داخل $(1 + \epsilon)Meb(K_1)$ قرار

بگیرند، ادامه می‌دهد. اولین نقطه‌ای که در شرایط ذکر شده صدق نمی‌کند را p_2 بنامید. حال الگوریتم کلارکسون را بر روی $K_1 \cup \{p_2\}$ اجرا کرده و مجموع هسته‌ی به دست آمده را K_2 بنامید. با ادامه‌ی این روند، الگوریتم، دنباله‌ای از مجموعه هسته $\kappa = \{K_1, \dots, K_u\}$ نگه می‌دارد و زمانی که نقطه‌ی p_{u+1} پیدا شود که در هیچ کدام از $Meb(K_j)$ ($1 \leq j \leq u$) نباشد، الگوریتم کلارکسون را برای $\bigcup_{j=1}^u K_j \cup \{p_{u+1}\}$ اجرا نموده و مجموعه هسته‌ی به دست آمده را K_{u+1} می‌نامد. با توجه به نحوه‌ی ساخته شدن K_i ها، به راحتی می‌توان نشان داد رابطه‌ی زیر برقرار است:

$$P \subset \bigcup_{i=1}^u (1 + \epsilon) Meb(K_i)$$

حال در نهایت برای به دست آوردن جواب نهایی کافی است کوچک‌ترین کره‌ای که

$$\bigcup_{i=1}^u (1 + \epsilon) Meb(K_i)$$

را می‌پوشاند را به عنوان جواب محاسبه کنیم. چن و سایرین ثابت کرده‌اند که کره‌ی نهایی دارای شعاع حداکثر $1/22$ برابر شعاع بهینه است. برای مشاهده جزئیات بیشتر، به [۲۴، ۹] مراجعه کنید.

آگاروال نه تنها الگوریتمی ارائه داد که در نهایت، ثابت شد حداکثر جوابی با ضریب تقریب $1/22$ برابر جواب بهینه می‌دهد، بلکه نشان داد، که با حافظه‌ی چندجمله‌ای بر اساس $\log n$ و d نمی‌توان الگوریتمی ارائه داد که ضریب تقریب بهتر از $\frac{1+\sqrt{2}}{4}$ داشته باشد.

قضیه ۳-۱ هر الگوریتم تحت مدل جویبار داده که یک α -تقریب برای مسئله‌ی ۱-مرکز برای مجموعه‌ی S شامل n نقطه در فضای \mathbb{R}^d نگه می‌دارد، برای $\alpha \leq \frac{1+\sqrt{2}}{4}(1 - \frac{2}{d^{\frac{1}{3}}})$ با احتمال حداقل $\frac{2}{3}$ نیاز به $\Omega(\min\{n, e^{d^{\frac{1}{3}}}\})$ حافظه مصرف می‌کند.

اثبات. ایده‌ی اصلی اثبات بر اساس قضیه‌ی معروف آلیس و باب^{۱۲} در نظریه انتقال اطلاعات بنا شده است. برای خواندن اثبات این قضیه می‌توانید به مرجع [۹] مراجعه کنید.

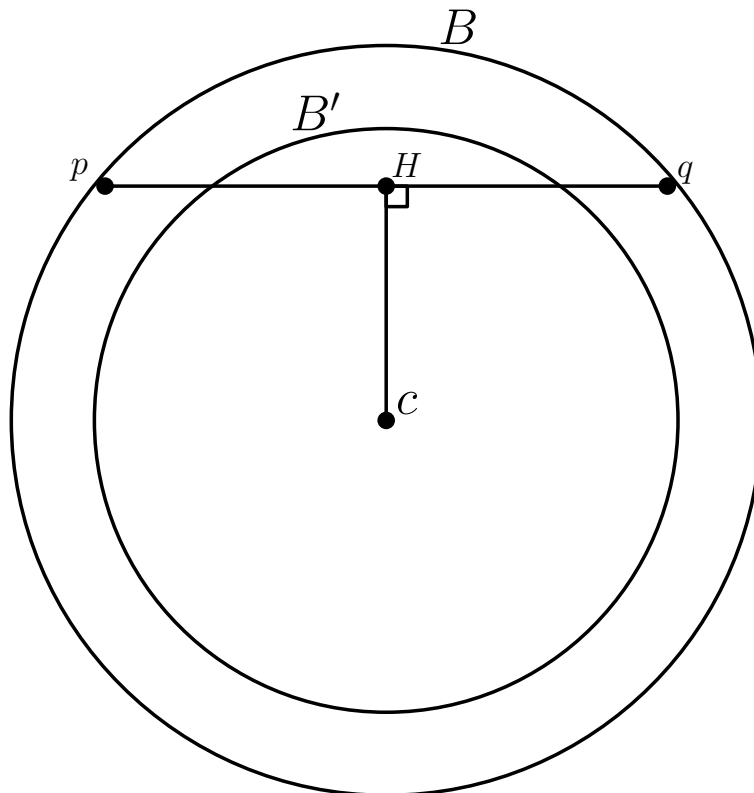
□

علاوه بر مسئله‌ی ۱-مرکز، مسئله‌ی دو مرکز نیز در سال‌های اخیر مورد توجه قرار گرفته است و بهبودهایی نیز برای این مسئله ارائه شده است. آهن و سایرین [۲۶] در سال ۲۰۱۴، اولین الگوریتم با

^{۱۲}alice and bob

ضریب تقریب کمتر از دو را برای مسئله ۲- مرکز در فضای اقلیدسی ارائه دادند. این الگوریتم تقریباً پایه‌ی کار این پایان‌نامه برای حالت‌های مختلف است. به همین منظور، تعدادی از لم‌های داخل این الگوریتم که در آینده استفاده می‌شود این‌جا توضیح داده می‌شود.

لم ۲-۳ فرض کنید B یک کره‌ی واحد با مرکز c در فضای اقلیدسی \mathbb{R}^d باشد. هر پاره‌خط pq به طول حداقل $1/2$ که به طور کامل داخل B قرار دارد، کره‌ی $B'(c, 1/8)$ را قطع می‌کند.



شکل ۳-۶: اثبات لم ۲-۳

اثبات. صفحه‌ی گذرنده از پاره‌خط و مرکز کره را نظر بگیرید. ادامه اثبات تنها به همین صفحه محدود می‌شود، بنابراین نیاز به در نظر گرفتن ابعاد بزرگ‌تر از ۲ نیست. همان‌طور که در شکل ۳-۶ مشخص شده است، پای عمود از مرکز کره بر پاره‌خط pq را h بنامید. بدون کم شدن از کلیت مسئله فرض می‌کنیم $\|hq\| \leq \|hp\|$. بنابراین داریم:

$$1/6 = \frac{1/2}{2} \leq \|hq\|$$

از طرفی چون پاره‌خط pq به طور کامل داخل کره‌ی واحد قرار گرفته است، بنابراین تمام نقاط آن، شامل

دو سر آن، از مرکز کره، فاصله‌ی حداکثر ۱ دارند. بنابراین طبق رابطه‌ی فیثاغورث، داریم:

$$\|hc\| = \sqrt{\|qc\|^2 - \|qh\|^2} \leq \sqrt{1 - 0/6^2} = 0/8$$

بنابراین نقطه‌ی h داخل کره‌ی B' قرار می‌گیرد. از طرفی چون، $\|pq\| \leq 1$ بنابراین، h داخل پاره‌خط قرار دارد و در نتیجه کره‌ی B' با پاره‌خط pq تقاطع دارد.

□

لم بالا در واقع نشان می‌دهد اگر در طول الگوریتم بتوانیم دو نقطه‌ی دور نسبت به هم (حداقل $1/2$ برابر شعاع بهینه) از یکی از دو کره‌ی بهینه را بیابیم، این پاره‌خط از مرکز کره‌ی بهینه فاصله‌ی کمی (حداکثر $0/8$ شعاع بهینه) دارد.

لم ۳-۳ فرض کنید B کره‌ای به مرکز c و شعاع واحد در \mathbb{R}^d باشد. پاره‌خط دلخواه pq با طول حداقل $1/2$ که به طور کامل داخل B قرار دارد را در نظر بگیرید. هر نقطه‌ی x از پاره‌خط pq که از دو سر آن حداقل $0/6$ فاصله داشته باشد، داخل کره‌ی $B'(c, 0/8)$ قرار می‌گیرد.

اثبات. اثبات این لم نیز کاملاً مشابه لم ۲-۳ است. بدون کم شدن از کلیت مسئله همان‌طور که در شکل ۷-۳ مشخص شده است، فرض کنید زاویه‌ی $\angle pxc$ بزرگ‌تر مساوی 90° درجه است. در نتیجه داریم:

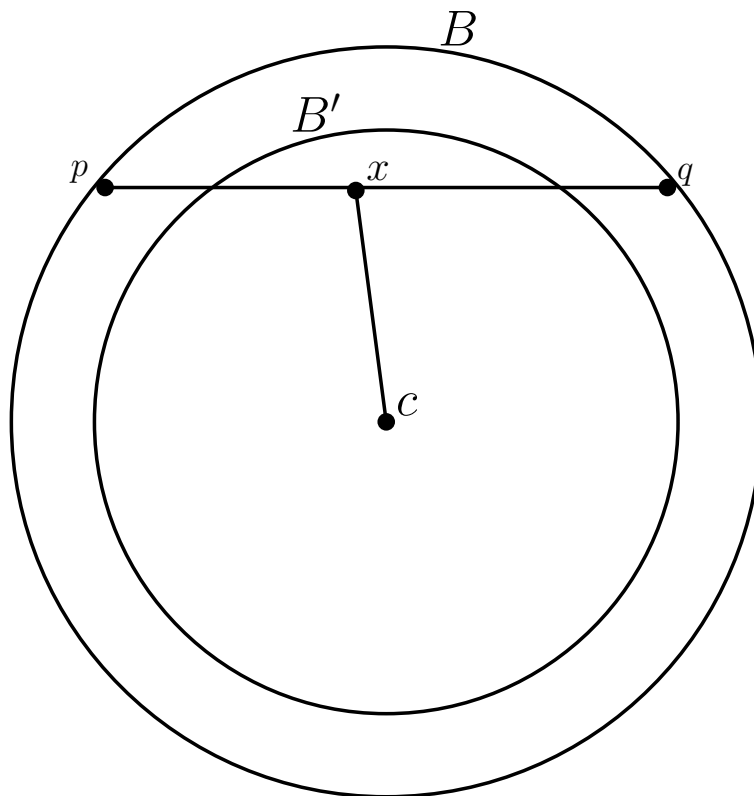
$$\sqrt{\|px\|^2 + \|xc\|^2} \leq \|pc\| \leq 1$$

از طرفی طبق فرض مسئله داریم:

$$0/6 \leq \|px\| \implies \|xc\| \leq \sqrt{1 - 0/6^2} = 0/8$$

□

الگوریتم آهن، با استفاده از دو لم بالا و تقسیم مسئله به دو حالتی که دو کره‌ی بهینه بیش از دو برابر شعاع بهینه یا کم‌تر از برابر شعاع بهینه فاصله داشته باشند، دو الگوریتم کاملاً جداگانه ارائه می‌دهد که به طور موازی اجرا می‌گردند. اجرای موازی این دو الگوریتم، در هر لحظه دو جواب درست ارائه می‌دهد و کافی است برای جواب نهایی بین دو شعاعی که به عنوان جواب خود می‌دهند، شعاع کم‌تر را به عنوان جواب نهایی الگوریتم بدهیم.

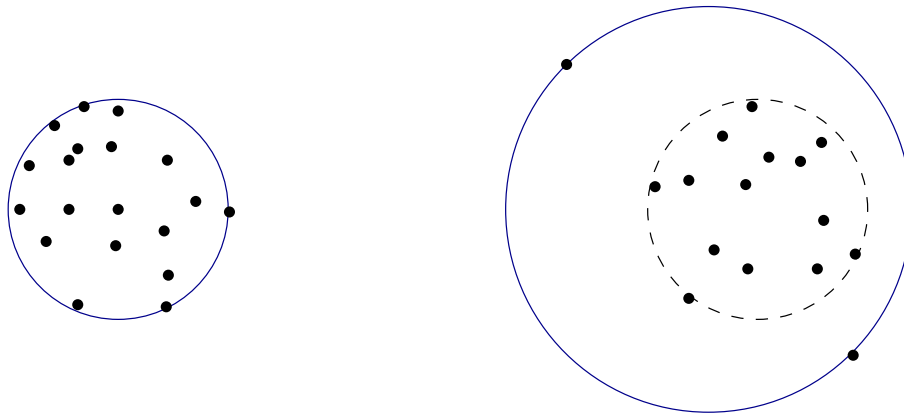


شکل ۳-۷: اثبات لم ۳-۳

۳-۳ - مرکز با داده‌های پرت k

در دنیای واقعی در میان داده‌ها، داده‌های دارای اریب وجود دارد که اگر امکان تشخیص و حذف آن‌ها در حین جمع‌آوری داده‌ها وجود داشت، شعاع مسئله‌ی k - مرکز به میزان قابل توجهی کاهش پیدا می‌کرد و شهود بسیار بهتری از دسته‌ها ارائه می‌داد. در عمل نیز حذف داده‌هایی که به هیچ دسته‌ای شباهت ندارند، معقول به نظر می‌رسد، زیرا هدف به دست آوردن دسته‌بندی برای کلیت نقاط است و نه دسته‌بندی که تمام نقاط در آن می‌گنجند. همان‌طور که در شکل ۳-۸ می‌بینید، تنها حذف دو نقطه که نسبت به بقیه نقاط داده‌ی اریب حساب می‌شوند، دسته‌بندی بسیار قابل قبول‌تری را نتیجه می‌دهد. با چنین رویکردی، باید تعدادی نقطه که با بقیه نقاط فاصله‌ی زیادی دارند از داده‌های ورودی حذف شده و سپس به عنوان ورودی مسئله‌ی k - مرکز دسته‌های مرتب را از آن استخراج کرد.

مسئله‌ی k - مرکز با داده‌های پرت، بسیار مشابه مسئله‌ی استقرار تجهیزات است. در مسئله‌ی استقرار تجهیزات، هدف استقرار چند مرکز ارائه دهنده‌ی خدمات است که هزینه‌ی استقرار به علاوه‌ی انتقال



شکل ۳-۸: کاهش قابل توجه شعاع مسئله ۲- مرکز با حذف تنها دو نقطه

تجهیزات از مراکز ارائه دهنده به مکان‌ها متقاضی کمینه گردد. تعریف رسمی این مسئله در زیر آمده است:

مسئله ۳-۱ (استقرار تجهیزات) مجموعه‌ای نقطه به عنوان مکان‌های مجاز برای استقرار تجهیزات داده شده است. هزینه استقرار تجهیزات در نقطه‌ی i ام را برابر با f_i در نظر بگیرید. مجموعه‌ای از نقاط نیز که متقاضی تجهیزات هستند نیز داده شده است. به ازای هر متقاضی j و محل استقرار تجهیزات i ، $d(i, j)$ برابر هزینه انتقال تجهیزات از محل استقرار به متقاضی است. مجموعه‌ای k عضوی به نام K انتخاب کنید به طوری که هزینه کلی را کمینه نماید:

$$\sum_{i \in K} f_i + \sum_{\text{all customers}} \min_{i \in K} d(i, j)$$

گونه‌های مختلفی از مسئله استقرار تجهیزات تعریف شده است. از جمله‌ی آن می‌توان به گونه‌های زیر اشاره نمود:

- هر مرکز ارائه‌دهنده حداکثر به تعداد مشخصی از متقاضیان می‌تواند تجهیزات انتقال دهد. در واقع در این روش سعی در استقرار تجهیزات به صورت متوازن است و متناسب با قدرت ارائه‌ی تجهیزات آن‌هاست.
- هزینه استقرار مرکز ارائه‌دهنده متناسب با تعداد متقاضیانی که پاسخ می‌دهد است. در این روش، معمولاً هر چه تعداد تقاضاهای یک مرکز بالاتر رود هزینه استقرار یا ساخت آن برای تامین چنین میزان درخواستی بالاتر می‌رود.

• هر متقاضی میزانی جریمه بابت عدم دریافت تقاضای خود مشخص کند. در این روش، هر متقاضی در صورت عدم دریافت تجهیزات مورد نیاز، میزانی جریمه مطالبه می‌کند و هدف کاهش مجموع هزینه‌ها به علاوه‌ی هزینه‌های قبلی است.

• تعدادی از متقاضیان را می‌توان بدون پرداخت جریمه پوشش نداد. در این حالت، امکان چشم‌پوشی از تعدادی از متقاضیان وجود دارد ولی امکان تخطی از محدودیت تعداد آن‌ها وجود ندارد.

اگر به دو گونه‌ی آخر توجه بیش‌تری کنید، به شباهتشان به مسئله‌ی k -مرکز با داده‌های پرت پی خواهید برد. می‌توان نشان داد که مسئله‌ی k -مرکز با z داده‌ی پرت از لحاظ پیچیدگی محاسباتی هم‌ارز استقرار تجهیزات با امکان عدم پوشش z متقاضی است. همان‌طور که در زیربخش قبلی دیدیم، هیچ الگوریتمی با ضریب تقریب کمتر از ۲ برای مسئله‌ی k -مرکز در حالت کلی وجود ندارد مگر این‌که $P = NP$ باشد. برای مسئله‌ی k -مرکز با داده‌های پرت، اگر تعداد مجاز داده‌های پرت صفر باشد، مسئله به همان مسئله‌ی k -مرکز تبدیل می‌گردد. گونه‌ی مشابهی با مسئله‌ی k -مرکز با داده‌های پرت تعریف می‌گردد که در آن تعدادی از نقاط نمی‌توانند به عنوان مرکز انتخاب شوند. به این‌گونه، مسئله‌ی k -مرکز با داده‌های پرت و داده‌های ممنوعه (برای قرارگیری مرکز در آن‌ها) تعریف می‌شود. در مرجع [۲]، ثابت شده است که این مسئله در حالت کلی با ضریب کمتر از ۳ قابل تقریب‌پذیر نیست مگر آن‌که $P = NP$ باشد.

قضیه‌ی ۳-۴ فرض کنید نقطه‌هایی دلخواه از یک متریک دلخواه داده شده‌اند. مسئله‌ی k -مرکز با داده‌های پرت، که در آن، بعضی از رئوس امکان مرکز شدن ندارند (رئوس ممنوعه)، را نمی‌توان با ضریب تقریبی کمتر از ۳ تقریب زد.

اثبات. ایده‌ی ارائه شده برای این کران پایین، بسیار مشابه با ایده‌ی ارائه شده برای مسئله‌ی k -مرکز در فضای اقلیدسی است [۸]. در واقع در این راه‌حل، مسئله‌ی بیشینه پوشش رأسی به یک گراف دوبخشی متریک تبدیل می‌گردد که در آن وزن تمام یال‌ها برابر یک است. با استفاده از گراف ساخته شده، نشان داده می‌شود که اگر مسئله‌ی k -مرکز با داده‌های پرت و مراکز ممنوعه، دارای الگوریتمی با ضریب تقریب کمتر از ۳ داشته باشد (کوتاه‌ترین مسیر بین دو رأس غیر مجاور در دو بخش مختلف)، آنگاه می‌توان مسئله‌ی بیشینه پوشش رأسی را در زمان چندجمله‌ای حل نمود. برای مشاهده‌ی اثبات کامل‌تر می‌توانید به مرجع [۳] مراجعه کنید.

□

الگوریتمی که چریکار^{۱۳} سایرین در این مقاله ارائه داده‌اند یک الگوریتم ۳- تقریب برای مسئله‌ی k -مرکز با داده‌های پرت است. در ابتدا، الگوریتم چریکار، تمام شعاع‌های ممکن را پیدا می‌کند. در حالت مسئله‌ی k -مرکز گسته، کافی است تمام فاصله‌های بین دویه‌دوی نقاط را به عنوان کاندیدا در نظر گرفت که تعدادشان از مرتبه‌ی $O(n^2)$ است و کافی است بر روی گزینه‌های به دست آمده، جست‌وجوی دودویی زد. در صورتی که مسئله‌ی k -مرکز در حالت پیوسته مد نظر باشد، به ازای هر $d + 1$ نقطه‌ی دلخواه، یک کاندید اضافه می‌گردد که تعدادشان از مرتبه‌ی $O(n^{d+1})$ است. حال اگر ما شعاع r داشته باشیم که r بزرگ‌تر مساوی شعاع بینه باشد، چریکار یک الگوریتم ساده با شعاع حداکثر $3r$ ارائه می‌دهد و اگر الگوریتم چریکار نتوانست با شعاع $3r$ همه‌ی نقاط به جز حداکثر z نقطه را بپوشاند، بنابراین فرض اولیه‌ی ما اشتباه بوده است و r از شعاع بهینه کم‌تر است.

تعریف ۲-۳ به ازای هر نقطه‌ی $v_i \in V$ ، G_i (به طور مشابه) را برابر مجموعه نقاطی در نظر بگیرید که در فاصله‌ی حداکثر r (به طور مشابه) از v_i قرار دارند. G_i را توپ به شعاع r و E_i را به عنوان توپ گسترش‌یافته به شعاع $3r$ می‌نامیم. وزن هر توپ را برابر تعداد نقاط درون آن در نظر می‌گیریم.

الگوریتم ۳ اولین الگوریتم با ضریب تقریب ۳ برای k -مرکز با z داده‌ی پرت

- ۱: به ازای i از ۱ تا k :
 - ۲: به ازای تمام نقاط، توپ‌ها و توپ‌های گسترش‌یافته را محاسبه کن.
 - ۳: G_j را توپی در نظر بگیر که بیش‌ترین وزن را دارد (بیش‌ترین تعداد نقاط پوشش داده نشده را می‌پوشاند).
 - ۴: توپ E_j را به عنوان توپ i ام در نظر بگیر.
 - ۵: تمام نقاط داخل E_j را به مجموعه نقاط پوشش داده شده اضافه کن.
 - ۶: اگر همه‌ی نقاط به جز حداکثر z تای آن‌ها پوشانده شده بودند:
 - ۷: r اولیه بزرگ‌تر مساوی r بهینه است. برگردان E_j ‌های انتخاب شده
 - ۸: در غیر این صورت:
 - ۹: r اولیه کوچک‌تر از r بهینه است.
-

الگوریتم ۳، یک الگوریتم حریصانه و ساده است که با مقایسه‌ی عمل‌کرد آن با جواب بهینه می‌توان نشان داد که به درستی عمل می‌کند. برای مشاهده‌ی درستی اثبات، می‌توانید به [۳] کنید. چریکار در ادامه، با استفاده از برنامه‌ریزی خطی^{۱۴} و گردسازی^{۱۵} جواب، یک الگوریتم ۳- تقریب برای حالت گسسته و یک ۴- تقریب برای حالت پیوسته ارائه می‌دهد. به علت عدم استفاده از روش برنامه‌ریزی خطی، از بیان جزئیات این قسمت صرف نظر می‌کنیم. مک‌کاتن^{۱۶} با استفاده از الگوریتم ارائه شده در بالا، یک الگوریتم جویبار داده ارائه داد که متناسب با اینکه از کدام الگوریتم استفاده کند، همان ضریب تقریب را برای حالت جویبار داده می‌دهد. ایده‌ی اصلی به کار رفته در این تبدیل، پردازش نقاط به صورت دسته‌های $O(kz)$ تایی و در نظر گرفتن نقاط آزاد به عنوان نقاطی هنوز مطمئن نیستیم نقطه‌ی پرت هستند و نقاطی که مطمئن هستیم باید پوشانده شوند. این الگوریتم حافظه‌ای از مرتبه‌ی $O(\frac{kz}{\epsilon})$ مصرف می‌کند. برای مشاهده جزئیات بیشتر به مرجع [۴] مراجعه کنید.

تا به اینجا مسئله‌ی k - مرکز را در حالت کلی بررسی کردیم. مسئله‌ی ۱- مرکز با داده‌های پرت به صورت جداگانه به وسیله‌ی ضربی‌زاده و سایرین مورد بررسی قرار گرفته است [۲۷]. در الگوریتم ارائه شده برای حالتی که $z = 1$ است، ضربی‌زاده یک الگوریتم با ضریب تقریب $1/48 \leq 1/22 \times \frac{\sqrt{2}+1}{4} \leq$ با حافظه‌ی مصرفی $O(\frac{d}{\epsilon^3} \log \frac{1}{\epsilon})$ ارائه داده است. به ازای z های کلی، الگوریتم دیگری با حافظه‌ی مصرفی $O(d^3 z + \frac{d}{\epsilon^3} \log \frac{1}{\epsilon})$ و ضریب تقریب $1/73 \leq 1/22 \times \sqrt{2} \leq$ ارائه شده است.

ایده‌ی اصلی که در این مقاله ارائه شده است، ارائه‌ی یک راه‌کار کلی برای مسائل جویبار داده است. در این راه‌کار، یک حافظه‌ی میان‌گیر^{۱۷} تعریف می‌شود که هر نقطه از جویبار داده به محض ورود به آن اضافه می‌شود. در صورتی که حافظه‌ی میان‌گیر، پر گردد، یکی از نقاط از حافظه استخراج شده و به الگوریتم زیرین داده می‌شود. الگوریتم زیرین یک الگوریتم جویبار داده برای مسئله‌ی ۱- مرکز بدون داده‌های پرت است. همان‌طور که در زیربخش مسئله‌ی k - مرکز در حالت جویبار داده بیان شد، بهترین الگوریتم موجود یک الگوریتم با حافظه‌ی مصرفی $O(d^3 z + \frac{d}{\epsilon^3} \log \frac{1}{\epsilon})$ و ضریب تقریب $1/22$ است.

عامل ثانویه‌ای که در ضریب تقریب نهایی الگوریتم تاثیر به سزایی دارد نحوه‌ی استخراج نقطه از حافظه‌ی میان‌گیر است. فرض کنید O_x مجموعه نقاطی از P (جویبار داده) باشند که در جواب بهینه به عنوان داده‌ی پرت انتخاب شدند و O مجموعه نقاطی که به اشتباه از حافظه‌ی میان‌گیر استخراج شدند

^{۱۴}Linear Programming^{۱۵}Rounding^{۱۶}McCutchen^{۱۷}Buffer

باشد، اگر داشته باشیم که

$$Meb((P - O_x) \cup O) \leq \beta Meb(P - O_x)$$

در نتیجه ضریب تقریب نهایی برابر $1/22\beta$ خواهد بود. نکته‌ی قابل توجه در اینجا است که طول حافظه‌ی میان‌گیر، در ضریب β تأثیر به‌سزایی دارد.

در حالت کلی z ، ایده‌ی اصلی برای استخراج یک نقطه از حافظه‌ی میان‌گیر، نقطه‌ی مرکزی^{۱۸} نقاط داخل حافظه‌ی میان‌گیر است. در واقع نزدیک‌ترین نقطه به نقطه‌ی مرکزی نقاط داخل حافظه‌ی میان‌گیر، استخراج می‌گردد و ثابت می‌شود با این شیوه‌ی استخراج $\beta \leq \sqrt{2}$ خواهد بود. برا مشاهده‌ی اثبات و جزئیات بیشتر به مرجع [۲۷] مراجعه کنید.

در فصل آتی، در ابتدا به پیشرفت‌هایی که برای مسئله‌ی ۱ – مرکز در حالت جویبار داده با داده‌های پرت در این پایان‌نامه ارائه شده است، خواهیم پرداخت. سپس برای مسئله‌ی ۲ – مرکز در حالت جویبار داده با داده‌های پرت، اولین کار موجود را ارائه می‌دهیم که بهبود قابل توجهی نسبت به حالت کلی است.

فصل ۴

نتایج جدید

در این فصل نتایج جدید به دست آمده در پایان نامه توضیح داده می شود. در صورت نیاز می توان نتایج جدید را در قالب چند فصل ارائه نمود. همچنین در صورت وجود پیاده سازی، بهتر است نتایج پیاده سازی را در فصل مستقلی پس از این فصل قرار داد.

فصل ۵

نتیجه‌گیری

در این فصل، ضمن جمع‌بندی نتایج جدید ارائه‌شده در پایان‌نامه، مسائل باز باقی‌مانده و همچنین پیشنهادهایی برای ادامه‌ی کار ارائه می‌شوند.

پیوست آ

مطالب تکمیلی

پیوست‌های خود را در صورت وجود می‌توانید در این قسمت قرار دهید.

کتاب نامه

- [1] *J. Han and M. Kamber. Data Mining, Southeast Asia Edition: Concepts and Techniques. Morgan kaufmann, 2006.*
- [2] *C. C. Aggarwal. Data streams: models and algorithms, volume 31. Springer Science & Business Media, 2007.*
- [3] *M. Charikar, S. Khuller, D. M. Mount, and G. Narasimhan. Algorithms for facility location problems with outliers. In Proceedings of the 12th ACM-SIAM Symposium on Discrete Algorithms, pages 642–651, 2001.*
- [4] *R. M. McCutchen and S. Khuller. Streaming algorithms for k -center clustering with outliers and with anonymity. In Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques, pages 165–178. 2008.*
- [5] *M. Sipser. Introduction to the Theory of Computation. Cengage Learning, 2012.*
- [6] *I. Dinur and S. Safra. On the hardness of approximating minimum vertex cover. Annals of mathematics, pages 439–485, 2005.*
- [7] *V. V. Vazirani. Approximation algorithms. Springer Science & Business Media, 2013.*
- [8] *M. Bern and D. Eppstein. Approximation algorithms for geometric problems. Approximation algorithms for NP-hard problems, pages 296–345, 1996.*
- [9] *P. K. Agarwal and R. Sharathkumar. Streaming algorithms for extent problems in high dimensions. In Proceedings of the 21st ACM-SIAM Symposium on Discrete Algorithms, pages 1481–1489, 2010.*
- [10] *P. K. Agarwal, S. Har-Peled, and K. R. Varadarajan. Approximating extent measures of points. Journal of the ACM, 51(4):606–635, 2004.*

- [11] R. G. Michael and S. J. David. *Computers and intractability: a guide to the theory of np-completeness*. WH Freeman & Co., San Francisco, 1979.
- [12] N. Megiddo and K. J. Supowit. *On the complexity of some common geometric location problems*. SIAM Journal on Computing, 13(1):182–196, 1984.
- [13] J. Han, M. Kamber, and J. Pei. *Data mining: concepts and techniques: concepts and techniques*. Elsevier, 2011.
- [14] P. K. Agarwal and C. M. Procopiuc. *Exact and approximation algorithms for clustering*. Algorithmica, 33(2):201–226, 2002.
- [15] N. Megiddo. *On the complexity of some geometric problems in unbounded dimension*. Journal of Symbolic Computation, 10(3):327–334, 1990.
- [16] B. Chazelle and J. Matoušek. *On linear-time deterministic algorithms for optimization problems in fixed dimension*. Journal of Algorithms, 21(3):579–597, 1996.
- [17] T. M. Chan. *More planar two-center algorithms*. Computational Geometry: Theory and Applications, 13(3):189–198, 1999.
- [18] P. K. Agarwal, R. B. Avraham, and M. Sharir. *The 2-center problem in three dimensions*. Computational Geometry, 46(6):734–746, 2013.
- [19] H. Zarrabi-Zadeh. *Core-preserving algorithms*. In Proceedings of the 20th Canadian Conference on Computational Geometry, pages 159–162, 2008.
- [20] M. Charikar, C. Chekuri, T. Feder, and R. Motwani. *Incremental clustering and dynamic information retrieval*. In Proceedings of the 29th ACM Symposium on Theory of Computing, pages 626–635. ACM, 1997.
- [21] S. Guha. *Tight results for clustering and summarizing data streams*. In Proceedings of the 12th International Conference on Database Theory, pages 268–275, 2009.
- [22] H.-K. Ahn, H.-S. Kim, S.-S. Kim, and W. Son. *Computing k centers over streaming data for small k*. International Journal of Computational Geometry and Applications, 24(02):107–123, 2014.
- [23] H. Zarrabi-Zadeh and T. M. Chan. *A simple streaming algorithm for minimum enclosing balls*. In Proceedings of the 18th Canadian Conference on Computational Geometry, pages 139–142, 2006.

-
- [24] T. M. Chan and V. Pathak. *Streaming and dynamic algorithms for minimum enclosing balls in high dimensions*. Computational Geometry: Theory and Applications, 47(2):240–247, 2014.
- [25] M. Badoiu and K. L. Clarkson. *Smaller core-sets for balls*. In Proceedings of the 14th ACM-SIAM Symposium on Discrete Algorithms, pages 801–802. Society for Industrial and Applied Mathematics, 2003.
- [26] S.-S. Kim and H.-K. Ahn. *An improved data stream algorithm for clustering*. In Proceedings of the 11th, pages 273–284. 2014.
- [27] H. Zarrabi-Zadeh and A. Mukhopadhyay. *Streaming 1-center with outliers in high dimensions*. In Proceedings of the 21st Canadian Conference on Computational Geometry, pages 83–86, 2009.

واژه‌نامه

الف

heuristic..... ابتکاری

worth..... ارزش

Abstract

This part should be completed.

Keywords: *Clustering, K-Center, Streaming Data, Approximation Algorithm*



Sharif University of Technology

Department of Computer Engineering

M.Sc. Thesis

***Approximation Algorithms for Clustering Points
in the Data Stream Model***

By:

Behnam Hatami-Varzaneh

Supervisor:

Dr. Hamid Zarrabi-zade

September 2015