



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
۱۳	۳-۱ ادبیات موضوع
۱۴	۴-۱ اهداف تحقیق
۱۵	۵-۱ ساختار پایان نامه
۱۶	۲ مفاهیم اولیه
۱۶	۱-۲ مسائل ان پی- سخت
۱۸	۲-۱-۱ پوشش رأسی
۱۹	۲-۲ الگوریتم های تقریبی
۲۰	۲-۲-۱ میزان تقریب پذیری مسائل
۲۱	۳-۲ الگوریتم های جویبار داده
۲۲	۳ کارهای پیشین
۲۲	۱-۳ k -مرکز در حالت ایستا
۲۵	۲-۳ k -مرکز در حالت جویبار داده

۳۲	۳-۳ k -مرکز با داده‌های پرت
۳۳		۴ نتایج جدید
۳۴		۵ نتیجه‌گیری
۳۵		آ مطالب تکمیلی

فهرست شکل‌ها

۱-۱	نمونه‌ای از مسئله‌ی ۲- مرکز	۱۱
۲-۱	نمونه‌ای از مسئله‌ی ۲- مرکز با داده‌های پرت	۱۱
۳-۱	نمونه‌ای از مسئله‌ی ۲- مرکز در حالت پیوسته	۱۲
۱-۳	نمونه‌ای از تخصیص نقاط به ازای مراکز آبی رنگ.	۲۳
۲-۳	نمونه‌ای از حل مسئله‌ی ۳- مرکز با الگوریتم گنزالز	۲۴
۳-۳	نمونه‌ای از توری‌بندی الگوریتم ضربابی‌زاده (نقاط آبی، مراکز به دست آمده از الگوریتم تقریبی است). پس از توری‌بندی کافی است برای هر کدام از خانه‌های شبکه‌بندی، تنها یک نقطه را در نظر بگیریم.	۲۶
۴-۳	نمونه‌ای از اجرای الگوریتم ضربابی‌زاده 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.	۲۷
۵-۳	اثبات لم ۲-۳	۳۰
۶-۳	اثبات لم ۳-۳	۳۱

فهرست جدول‌ها

۱-۲ نمونه‌هایی از ضرایب تقریب برای مسائل بهینه‌سازی ۲۰

فصل ۱

مقدمه

مسئله‌ی خوشه‌بندی یکی از مهم‌ترین مسائل داده‌کاوی^۱ 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^۲، پردازش تصویر^۳، بازیابی اطلاعات^۴ و بایوانفورماتیک^۵ مورد استفاده قرار می‌گیرد [۱].

مسئله‌ی خوشه‌بندی، به‌خودی‌خود یک مسئله‌ی الگوریتمی به حساب نمی‌آید. راه‌حل‌های الگوریتمی بسیار زیادی برای خوشه‌بندی تعریف شده است. به طور کلی این الگوریتم‌ها به چهار دسته‌ی زیر تقسیم‌بندی می‌گردند:

- خوشه‌بندی‌های سلسه‌مراتبی^۶
- خوشه‌بندی‌های مرکزگرا^۷
- خوشه‌بندی‌های مبتنی بر توزیع^۸ نقاط

^۱ data mining

^۲ pattern recognition

^۳ image analysis

^۴ information retrieval

^۵ bioinformatics

^۶ hierarchical clustering

^۷ centroid-based clustering

^۸ distribution-based

• خوشه‌بندی‌های مبتنی بر چگالی^۹ 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 خوشه‌بندی‌های 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^{۱۰} یا مسئله‌ی استقرار تجهیزات متریک^{۱۱} یک مسئله‌ی بهینه‌سازی ترکیبیاتی^{۱۲} است. فرض کنید که n شهر و فاصله‌ی دوبه‌دوی 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^{۱۳} باشند و یا به زبان دیگر، یک گراف کامل داشته باشیم که فاصله‌ها در آن در رابطه‌ی مثلثی^{۱۴} 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.

۱-۱ تعریف مسئله

تعریف دقیق‌تر مسئله‌ی k -مرکز در زیر آمده است:

مسئله‌ی ۱-۱ یک گراف کامل بدون جهت $G = (V, E)$ با فاصله‌های $d(v_i, v_j)$ که از نامساوی مثلثی پیروی می‌کنند داده شده است. زیرمجموعه $S \subseteq V$ با اندازه‌ی k را انتخاب کنید به طوری که عبارت زیر

^۹ density-based

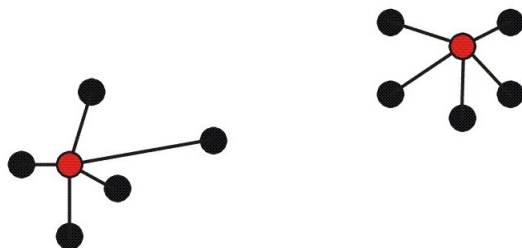
^{۱۰} metric

^{۱۱} metric facility location

^{۱۲} combinatorial optimization

^{۱۳} metric space

^{۱۴} triangle equation

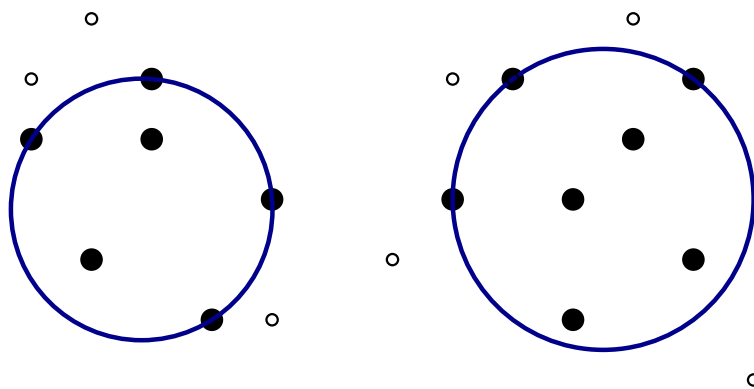


شکل ۱-۱: نمونه‌ای از مسئله‌ی ۲-مرکز

را کمینه کند:

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

گونه‌های مختلفی از مسئله‌ی k -مرکز با محدودیت‌های متفاوتی به‌وسیله‌ی پژوهشگران مورد مطالعه قرار گرفته است. از جمله می‌توان حالتی که در بین داده‌های ورودی داده‌های پرت وجود دارد و قبل از خوشه‌بندی می‌توانیم تعدادی نقاط ورودی را حذف نماییم و سپس به خوشه‌بندی بپردازیم. 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 به‌عنوان 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:



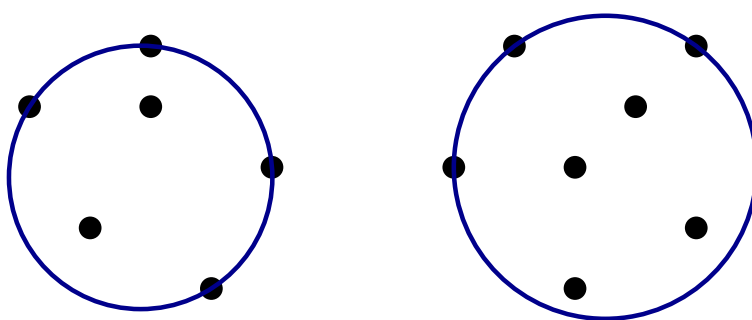
شکل ۱-۲: نمونه‌ای از مسئله‌ی ۲-مرکز با داده‌های پرت

مسئله ۱-۲ یک گراف کامل بدون جهت $G = (V, E)$ با فاصله‌های $d(v_i, v_j)$ که از نامساوی مثلثی پیروی می‌کنند داده شده است. زیرمجموعه $S \subseteq V$ با اندازه‌ی k و $Z \subseteq V$ به اندازه‌ی z را انتخاب کنید به طوری که عبارت زیر را کمینه کند:

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

گونه‌ای دیگری که در مسئله‌ی k -مرکز که در سال‌های اخیر مورد توجه قرار گرفته است، حالت جویبار داده‌ی آن است. در این گونه، در ابتدا تمام نقاط در دسترس نیست، 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 به طور معمول باید مرتبه‌ی حافظه‌ای کم‌تر از مرتبه حافظه‌ی خطی^{۱۵} متناسب با تعداد نقاط استفاده نمود. 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.

یکی از دغدغه‌هایی که در مسائل جویبار داده وجود دارد، عدم داشتن تمام نقاط است. بنابراین در بعضی موارد ممکن است برای مسئله‌ی k -مرکز، مرکزی برای یک دسته انتخاب شود که در بین نقاط ورودی نیست. نمونه‌ای از مسئله‌ی ۲-مرکز در حالت پیوسته را در شکل ۱-۳ نشان داده شده است. این مسئله تنها برای L_p -متریک مطرح می‌شود، زیرا مرکز دسته‌ها ممکن است در هر نقطه از فضا قرار بگیرد و ما نیاز داریم که فاصله‌ی آن را از تمام نقاط بدانیم. تعریف دقیق مسئله در زیر آمده است:



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

مسئله ۱-۳ مجموعه‌ی U از نقاط فضای d بعدی داده شده است. زیرمجموعه $S \subseteq U$ با اندازه‌ی k

^{۱۵}sublinear

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 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 گیرد.

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، مدل جویبار داده‌ی مسئله در سال‌های اخیر مورد توجه قرار گرفته است. 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 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 Kmart's اعلام کرده است که به ۸۸ درصد جمعیت 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 زیادی 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 مختلف مسئله‌ی k - مرکز هستند و در حالت‌های خاصی به مسئله‌ی k - مرکز کاهش پیدا می‌کنند. مسئله‌ی k - مرکز در

حوزه‌ی مسائل ان‌پی-سخت^{۱۶} قرار می‌گیرد و با فرض $P \neq NP$ الگوریتم دقیق با زمان چندجمله‌ای برای آن وجود ندارد. بنابراین برای حل کارای^{۱۷} این مسائل از الگوریتم‌های تقریبی^{۱۸} استفاده می‌شود. برای مسئله‌ی k -مرکز، دو الگوریتم تقریبی معروف وجود دارد. در الگوریتم اول، که به روش حریصانه^{۱۹} 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^{۲۰}، الگوریتمی با ضریب تقریب ۲ ارائه می‌گردد. همچنین ثابت شده است، که بهتر از این ضریب تقریب، الگوریتمی نمی‌توان ارائه شود 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 قرار دهد.

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 می‌پوشاند. 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

^{۱۶} NP-hard^{۱۷} efficient^{۱۸} Approximation Algorithm^{۱۹} greedy^{۲۰} dominating set

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 $\epsilon + 4$ است. ما با ارائه الگوریتمی جدید، الگوریتمی با ضریب تقریب $\epsilon + 1/8$ برای این مسئله ارائه شد که بهبودی قابل توجه برای این مسئله است.

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. 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 مورد بررسی قرار می‌گیرد. در بخش دوم، حالت جویبار داده‌ی مسئله و مجموعه هسته‌های^{۲۱} مطرح برای این مسئله مورد بررسی قرار می‌گیرد. در نهایت، در بخش سوم، مسئله‌ی k -مرکز با داده‌های پرت مورد بررسی قرار می‌گیرد.

در فصل چهارم، نتایج جدیدی که در این پایان‌نامه به دست آمده است، ارائه می‌شود. این نتایج شامل الگوریتم‌های جدید برای مسئله‌ی ۱- مرکز در حالت جویبار داده، مسئله‌ی ۱- مرکز با داده‌های پرت در حالت جویبار داده و مسئله‌ی دو مرکز در حالت جویبار داده با داده‌های پرت می‌شود.

در فصل پنجم، به جمع‌بندی کارهای انجام شده در این پژوهش و ارائه‌ی پیشنهادهایی برای انجام کارهای آتی و تعمیم‌هایی که از راه‌حل 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 به وسیله‌ی تورینگ^۳ ارائه شد، این مسئله معادل پیدا کردن الگوریتمی برای این مسئله است که اثبات می‌شود امکان‌پذیر نیست [۴]. برخلاف مثال بالا، عمده‌ی مسائل علوم کامپیوتر از نعد 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 ان پی (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 را الگوریتم ارائه دهد. 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 [۴] مراجعه نمود.

^۴Efficiency

^۵Complexity theory

^۶formal

^۷Decision problems

^۸verifiable

همان‌طور که می‌دانید درستی یا عدم درستی $P \subset NP$ از جمله معروف‌ترین 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 می‌شوند و در صورتی که زمانی، خلاف این فرض اثبات گردد، آنگاه قسمت عمده‌ای از علوم کامپیوتر زیر سوال می‌رود.

در نظریه پیچیدگی، برای دسته‌بندی مسائل، یکی از روش‌های دسته‌بندی کاهش چندجمله‌ای^{۱۰} مسائل به یک‌دیگر است.

تعریف ۱-۲ می‌گوییم مسئله‌ای A در زمان چندجمله‌ای به مسئله‌ی B کاهش می‌یابد، اگر وجود داشته باشد الگوریتم چندجمله‌ای C که به ازای هر ورودی α برای مسئله‌ی A ، یک ورودی β در زمان چندجمله‌ای برای مسئله‌ی B بسازد، به طوری که A ، α را می‌پذیرد 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 چندجمله‌ای از واژه‌ی کاهش استفاده می‌کنیم. در پی جست‌جوهای که برای برابری دسته‌ی پی و ان‌پی صورت گرفت، مجموعه‌ای از مسائل که عمدتاً داخل ان‌پی هستند استخراج گردید که اگر ثابت شود یکی از آن‌ها متعلق به پی است، آنگاه تمام مسائل دسته‌ی ان‌پی متعلق به پی خواهند بود و در نتیجه $P = NP$ خواهد بود. به این مجموعه مسائل ان‌پی-سخت می‌گویند. در واقع مسائل این دسته، مسائلی هستند که تمام مسائل داخل دسته‌ی ان‌پی، به آن‌ها کاهش می‌یابند.

کوک و لوین در قضیه‌ای به نام **قضیه‌ی کوک-لوین** ثابت کردند مسئله‌ی صدق‌پذیری^{۱۱} یک مسئله‌ی ان‌پی-سخت است [۴]. با پایه قرار دادن این اثبات و استفاده از تکنیک کاهش، اثبات ان‌پی-سخت بودن سایر مسائل، بسیار ساده‌تر گردید. در ادامه مسئله‌ی پوشش رأسی^{۱۲} را تعریف می‌کنیم.

۲-۱-۱ پوشش رأسی

در این پایان‌نامه، از این مسئله به عنوان مسئله‌ی پایه برای اثبات ان‌پی-سخت بودن مسئله‌ی k -مرکز استفاده می‌شود. تعریف این مسئله مطابق زیر است:

^۹Open problem

^{۱۰}Polynomial Reduction

^{۱۱}Satisfiability problem

^{۱۲}Vertex Coverage

تعریف ۲-۲ گراف بدون جهت $G(V, E)$ داده شده است. هدف مسئله پیدا کردن مجموعه‌ی $S \subset VS$ 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 یک از همسایگانش داخل مجموعه‌ی S قرار دارد. نسخه‌ی تصمیم‌گیری این مسئله به این گونه تعریف می‌شود که آیا گراف داده‌شده دارای پوشش رأسی با اندازه‌ی k است.

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. برای مطالعه‌ی روند اثبات به مرجع [۴] 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 کارامدی نیستند. در عمل، عمده‌ی مسائل کاربردی به این دسته تعلق می‌گیرند و هیچ راه‌حل یا الگوریتم کارامدی برایشان وجود ندارد. یکی از رویکردهای رایج در برابر چنین مسائلی، صرف نظر کردن از دقت راه‌حل‌هاست. به طور مثال راه‌حل‌های مکاشفه‌ای^{۱۳} گوناگونی برای مسائل مختلف 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 خوب عمل کردن دارند و در عمل معمولاً 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 مطرح شد، سعی در ارائه‌ی الگوریتم‌های مکاشفه‌ای شد که تضمین

^{۱۳} heuristic

مسئله	کران پایین تقریب پذیری
پوشش رأسی	$1,3606 [5]$
k - مرکز	$2 [6]$
۱ - مرکز در حالت جویبار داده	$\frac{1+\sqrt{2}}{4} [7]$
k - مرکز با نقاط پرت و نقاط اجباری	$3 [3]$

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

می‌کنند اختلاف زیادی با الگوریتمی که جواب بهینه می‌دهد، ندارند. در واقع این الگوریتم‌ها همواره و در هر شرایطی، تقریبی از جواب بهینه را ارائه می‌دهند. به چنین الگوریتم‌هایی، **الگوریتم‌های تقریبی**^{۱۴} می‌گویند. 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 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. معمولاً چنین الگوریتم‌هایی حافظه‌ی مصرفی از مرتبه‌ی زیرخطی^{۱۵} دارند و به همین دلیل برای داده‌های حجیم بسیار کاربرد دارند.

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 است استفاده از الگوریتم‌های تقریبی برای حل آن مسئله است. یکی از عمده‌ترین دغدغه‌های مطرح در الگوریتم‌های تقریبی کاهش ضریب تقریب است و یا حتی امکان ارائه‌ی الگوریتم تقریبی با ضریبی ثابت. به طور 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،

^{۱۴} Approximation Algorithm

^{۱۵} sublinear

معمولا می‌توان کران پایینی برای میزان تقریب‌پذیری آن‌ها ارائه داد. در واقع هر مسئله‌ی ان‌پی، علاوه بر این الگوریتم کارآمدی برای حل آن وجود ندارد، بعضا الگوریتم تقریبی برای حل آن با ضریبی تقریب کمی نیز وجود ندارد. در جدول ۱-۲ از میزان تقریبی مسائل مختلفی که در این پایان‌نامه مورد استفاده قرار می‌گیرد ببینید.

۲-۳ الگوریتم‌های جویبارداده

فصل ۳

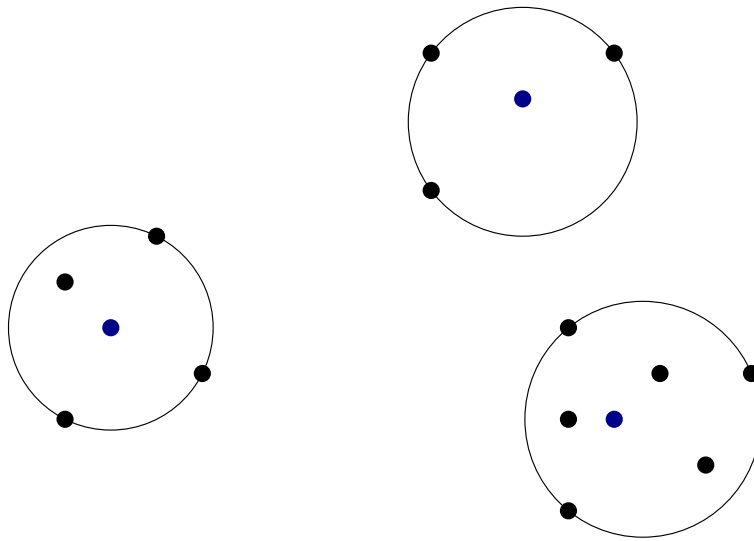
کارهای پیشین

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 بهینه سازی 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، اثبات گردید که این مسئله یک مسئله‌ی ان پی- سخت است [۸]. حتی ثابت شده است که این مسئله در صفحه‌ی دو بعدی و با متریک اقلیدسی نیز ان پی- سخت است [۹]. 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



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

است. برای چنین کاهشی کافی است، از روی گراف اصلی، یک گراف کامل بسازیم به طوری که معادل یال‌های گراف اصلی یال با وزن یک و به ازای بقیه یال‌های ممکن که در گراف اصلی نیستند، یال با وزن ۲ قرار می‌دهیم. 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 الگوریتم تقریبی با ضریب تقریب بهتر از $1/822$ وجود ندارد [۱۰].

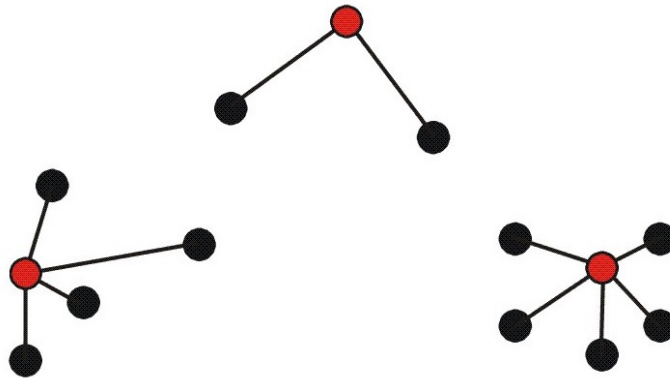
گنزالز^۲ اولین الگوریتم تقریبی برای مسئله‌ی k -مرکز را ارائه داده است [۱۱]. این الگوریتم یک الگوریتم تقریبی با ضریب ۲ است و در زمان $O(kn)$ قابل اجراست. الگوریتم گنزالز، یک الگوریتم حریصانه^۳ است. روش اجرای این الگوریتم به این گونه است که در ابتدا یک نقطه‌ی دلخواه را به عنوان مرکز دسته‌ی اول در نظر می‌گیرد. 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 می‌شود. 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

^۱ Euclidean space

^۲ Gonzalez

^۳ greedy

دسته، حداکثر دو برابر شعاع بهینه برای مسئله k -مرکز است. فدر^۴ و سایرین، زمان اجرای الگوریتم گنزالز را به $O(n \log k)$ برای هر L_p -متریک بهبود بخشیدند. نمونه‌ای از اجرای الگوریتم گنزالز، در شکل ۲-۳ نشان داده شده است.



شکل ۲-۳: نمونه‌ای از حل مسئله k -مرکز با الگوریتم گنزالز

تا به اینجا تنها برای k و d دلخواه صحبت شد. برای حالتی که در مسئله k -مرکز، k تعداد دسته‌ها و d ابعاد فضا ثابت باشند، 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 ارائه داده‌اند [۱۲]. قابل توجه است که اگر d ثابت نباشد، مسئله k -مرکز حتی برای متریک اقلیدسی (L_2 -متریک) با تعداد دسته‌ی ثابت $k \geq 2$ ، ان‌پی-سخت است [۱۳].

تا به اینجا، نتایج کلی برای k و d را مورد بررسی قرار دادیم. برای بعضی از مقادیر خاص از k و d الگوریتم‌های بهینه‌تری وجود دارد. به طور مثال، برای مسئله 1 -مرکز در فضای اقلیدسی با ابعاد ثابت، الگوریتم خطی با زمان اجرای $O((d+1)n)$ وجود دارد [۱۴]. الگوریتم ارائه شده بر پایه‌ی دو نکته‌ی اساسی بنا شده است. اول اینکه دایره‌ی 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 $O(\frac{d}{n})$ است که با توجه به ثابت بودن d این احتمال کوچک می‌باشد.

در صفحه اقلیدسی (L_2 -متریک) برای مسئله 2 -مرکز، بهترین الگوریتم را چن^۶ با زمان اجرای

^۴Feder

^۵Agarwal

^۶Chan

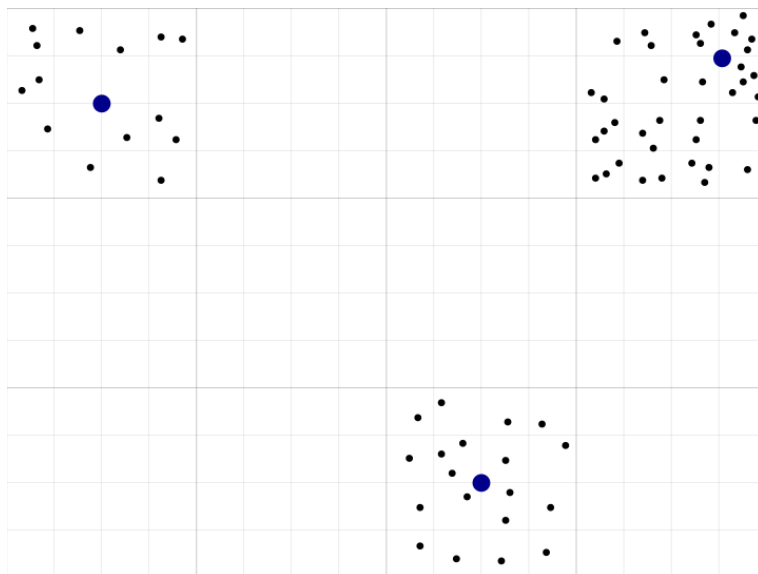
$O(c \log^2 n \log^2 \log n)$ و حافظه‌ی $O(n)$ ارائه داده است [۱۵]. برای فضای سه‌بعدی اقلیدسی نیز آگاروال و سایرین، الگوریتمی با متوسط زمان اجرای $O(n^3 \log^4 n)$ ارائه داده است [۱۶].

۳-۲ - مرکز در حالت جویبار داده k

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 ضرابی‌زاده برای نگه‌داری یک ϵ -هسته با حافظه‌ی $O(\frac{k}{\epsilon^d})$ برای L_p -متریک‌ها ارائه داده است [۱۷]. در روش ارائه شده، از چند ایده‌ی ترکیبی استفاده شده است. در ابتدا، الگوریتم با استفاده از الگوریتم تقریبی یک تقریب از جواب بهینه به دست می‌آورد. به طور مثال با استفاده از الگوریتم گزنالز، یک دو تقریب از شعاع بهینه به علاوه‌ی مرکز دسته‌های پیدا شده را به ما بدهد. حال کافی است که با طول شعاع الگوریتم دو تقریب حول هر مرکز، یک توری با شبکه‌بندی $O(\frac{1}{\epsilon})$ در هر بعد تشکیل دهیم و چون هر نقطه در حداقل یکی از توری‌ها قرار می‌گیرد، می‌توانیم 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 دوبرار سازی رایج در الگوریتم‌های جویبار داده است. نمونه‌ای از اجرای الگوریتم ضرابی‌زاده را در شکل ۳-۳ نشان داده شده است. برای دیدن اثبات‌ها و توضیح بیشتر در مورد روش ارائه شده می‌توانید به مرجع [۱۷] مراجعه کنید.

از جمله مشکلات وارده به الگوریتم ضرابی‌زاده، وابستگی سائز مجموعه‌ی هسته به ابعاد فضا است. بنابراین هسته‌ی ارائه شده به وسیله‌ی ضرابی‌زاده را نمی‌توان برای ابعاد بالا مورد استفاده قرار داد. از طرفی، حساب کردن جواب از روی هسته در زمان چندجمله‌ای قابل انجام نیست. با توجه با موارد گفته شده، برای قابل استفاده شدن الگوریتم‌ها برای ابعاد بالا، الگوریتم‌هایی ارائه می‌شود که ضریب تقریب



شکل ۳-۳: 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 یک نقطه را در نظر بگیریم.

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

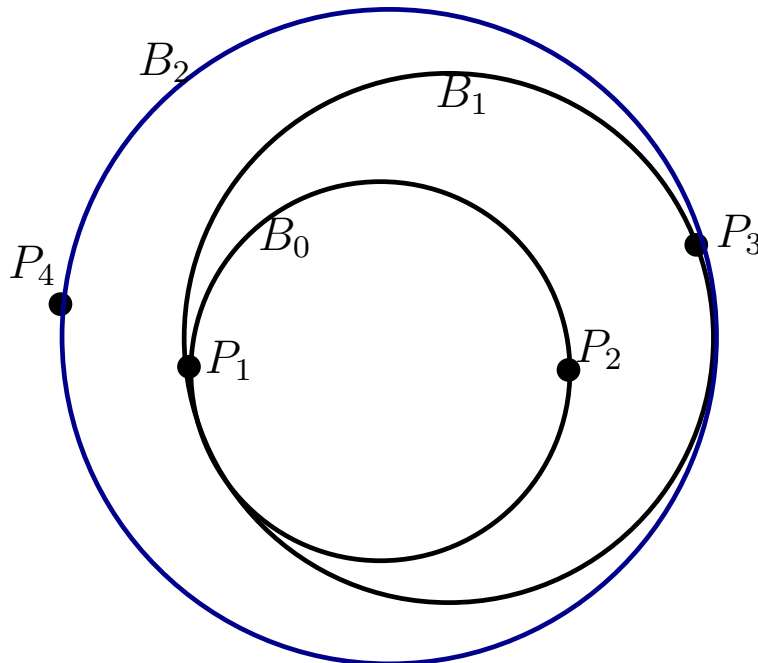
اولین الگوریتم ارائه شده برای ابعاد بالا، الگوریتمی با ضریب تقریب ۸ و حافظه‌ی مصرفی $O(dk)$ است [۱۸]. پس از آن، گوها^۷، به طور موازی با مک‌کاتن^۸ و سایرین، الگوریتمی با ضریب تقریب $(2 + \epsilon)$ با حافظه‌ی $O(\frac{dk}{\epsilon} \log \frac{1}{\epsilon})$ برای مسئله‌ی k -مرکز در هر فضای متریکی ارائه دادند [۱۹، ۲۰]. در سال ۲۰۱۴، اهن^۹ و سایرین، الگوریتمی با همین ضریب تقریب و حافظه‌ی $O((k+3)! 2^k \frac{d}{\epsilon})$ ارائه داده‌اند که برای k های ثابت، حافظه را از مرتبه‌ی $O(\log \frac{1}{\epsilon})$ کاهش می‌دهد [۲۱].

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

^۷Guha
^۸McCutchen
^۹Ahn

می‌توان یک هسته با اندازه $O(\frac{1}{\epsilon^{\frac{1}{d-1}}})$ ، با استفاده از نقاط حدی^{۱۰} در تعداد مناسبی جهت، به صورت جویبار داده ساخت.

مشکل عمده‌ی مجموعه هسته‌ی ارائه‌شده، وابستگی حافظه‌ی مصرفی آن به d است. ضربانی‌زاده و سایرین [۲۲] برای ابعاد دلخواه 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 می‌دارد، که کم‌ترین حافظه‌ی ممکن برای مسئله‌ی ۱-مرکز است. الگوریتم ارائه شده نقطه‌ی اول را به عنوان مرکز با شعاع صفر در نظر می‌گیرد. با فرا رسیدن هر نقطه‌ی جدید، اگر نقطه‌ی مد نظر، داخل دایره‌ی فعلی بیفتد که بدون هیچ تغییری ادامه می‌دهیم و در صورتی که بیرون دایره فعلی بیفتد، آن را با کوچک‌ترین دایره‌ای که نقطه‌ی جدید به علاوه‌ی دایره‌ی قبلی را به طور کامل می‌پوشاند، جایگزین می‌کنیم. به وضوح در هر لحظه دایره‌ی ساخته شده تمام نقاط را می‌پوشاند. از طرفی ثابت می‌شود شعاع دایره در هر لحظه، حداکثر $1/5$ - برابر شعاع دایره‌ی ۱-مرکز بهینه است. نمونه‌ای از اجرای الگوریتم را بر روی چهار نقطه می‌توان در شکل ۳-۴ دید. برای اثبات کامل‌تر می‌توانید به مرجع [۲۲] مراجعه کنید.



شکل ۳-۴: نمونه‌ای از اجرای الگوریتم ضربانی‌زاده 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 [۷] الگوریتمی تقریبی با حافظه‌ی مصرفی $O(d)$ ارائه دادند. در الگوریتم

^{۱۰}extreme points

ارائه شده، ضریب تقریب، برابر $\frac{1+\sqrt{3}}{4}$ تخمین زده شد، اما با تحلیل دقیق‌تری که چن و سایرین [۲۳] انجام دادند، مشخص شد، همان الگوریتم دارای ضریب تقریب $1/22$ است.

الگوریتم آن‌ها از الگوریتم کلارکسون و سایرین [۲۴] 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.

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 نقطه‌ی 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 می‌شود که در هیچ کدام از $(1+\epsilon)Meb(K_j)$ به ازای $1 \leq j \leq u$ 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} بنامید. 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 دایره‌ای که $(1+\epsilon)Meb(K_i)$ را می‌پوشاند را به عنوان جواب دهیم. آگاروال و سایرین ثابت کرده‌اند که دایره‌ی نهایی دارای شعاع حداکثر $1/22$ برابر شعاع بهینه است. برای مشاهده جزئیات بیش‌تر، به [۷] 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.

قضیه‌ی 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 $\alpha \leq \frac{1+\sqrt{2}}{4}(1 - \frac{2}{d^{\frac{1}{4}}})$ با احتمال حداقل $\frac{2}{3}$ نیاز به $\Omega(\min\{n, e^{d^{\frac{1}{4}}}\})$ حافظه مصرف می‌کند.

اثبات. ایده‌ی اصلی اثبات بر اساس قضیه‌ی معروف آلیس و باب^{۱۱} در نظریه انتقال اطلاعات بنا شده است. برای خواندن اثبات این قضیه می‌توانید به مرجع [۷] مراجعه کنید.

□

علاوه بر مسئله‌ی ۱- مرکز، مسئله‌ی دو مرکز نیز در سال‌های اخیر مورد توجه قرار گرفته است و بهبودهایی نیز برای آن مسئله ارائه شده است. آهن و سایرین [۲۵] در سال ۲۰۱۴، اولین الگوریتم با ضریب تقریب کمتر از دو را برای مسئله‌ی 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 $\frac{1}{2}$ که به طور کامل داخل B قرار دارد، دایره‌ی $B'(c, 1/8)$ را قطع می‌کند.

اثبات. صفحه‌ی گذرنده از پاره‌خط و مرکز دایره را نظر بگیرید. ادامه اثبات تنها به به همین صفحه محدود می‌شود، بنابراین نیاز به در نظر گرفتن ابعاد بزرگ‌تر از ۲ نیست. همان‌طور که در شکل ۳-۵ مشخص شده است، پای عمود از مرکز دایره بر پاره‌خط pq را h بنامید. بدون کم شدن از کلیت مسئله فرض می‌کنیم $\|hq\| \leq \|hp\|$. بنابراین داریم:

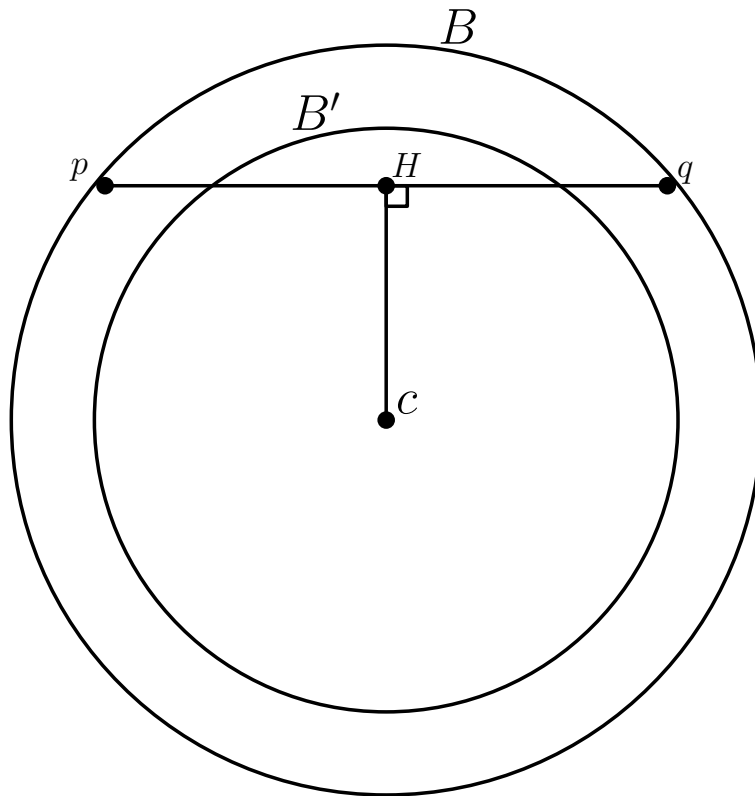
$$\frac{1}{6} = \frac{1/2}{4} \leq \|hq\|$$

از طرفی چون پاره‌خط pq به طور کامل داخل کره‌ی واحد قرار گرفته است، بنابراین تمام نقاط آن، شامل دو سر آن، از مرکز دایره، فاصله‌ی حداکثر ۱ دارند. بنابراین طبق رابطه‌ی فیثاغورث، داریم:

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

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

^{۱۱}alice and bob



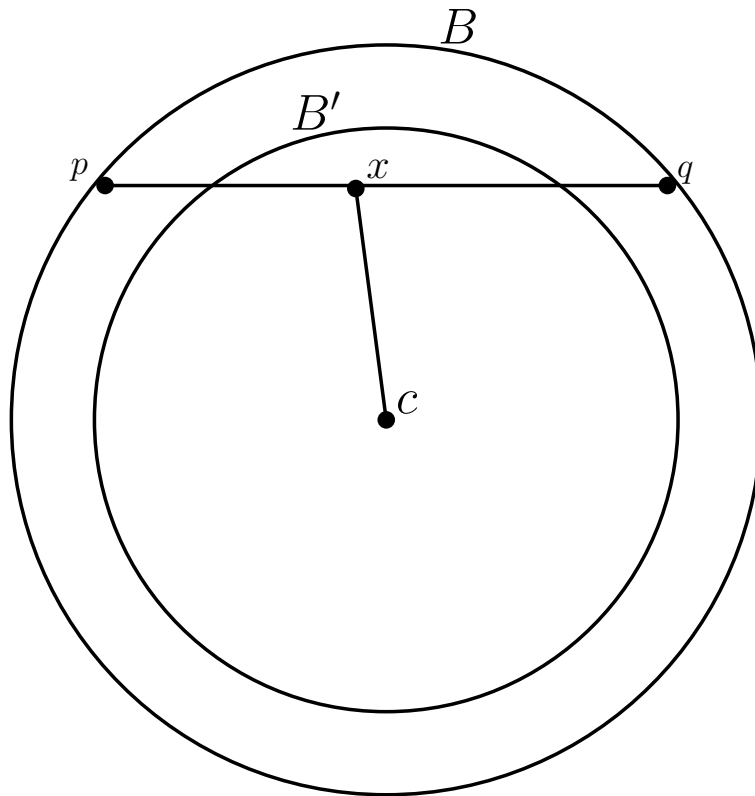
شکل ۳-۵: اثبات لم ۲-۳

□

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 دایره‌ی بهینه را بیاییم، این پاره‌خط از مرکز دایره‌ی بهینه فاصله‌ی کمی (حداکثر $0/8$ شعاع بهینه) را دارد.

لم ۳-۳ فرض کنید B دایره‌ای 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 دایره‌ی $B'(c, 0/8)$ قرار می‌گیرد.

اثبات. اثبات این لم نیز کاملاً مشابه لم ۲-۳ است. بدون کم شدن از کلیت مسئله همان‌طور که در شکل ۳-۶ مشخص شده است، فرض کنید زاویه‌ی $\angle pxc$ بزرگ‌تر مساوی 90° درجه است. در نتیجه



شکل ۳-۶: اثبات لم ۳-۳

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 دایره‌ی بهینه بیش از $2r^*$ یا کم‌تر از $2r^*$ فاصله داشته باشد، دو الگوریتم کاملاً جدا ارائه می‌دهد که به طور موازی اجرا می‌گردند. این دو الگوریتم، به طور موازی اجرا می‌شوند و در هر لحظه جوابی درست را به عنوان جواب نهایی ارائه می‌دهند و کافی است برای جواب نهایی بین دو شعاعی که به عنوان جواب خود می‌دهند، شعاع کم‌تر را به عنوان جواب الگوریتم بدهیم. به علت تشابه این الگوریتم، با یکی از الگوریتم‌های فصل کارهای جدید، از تکرار آن خودداری کرده و به تفصیل در فصل نتایج جدید شرح داده می‌شود.

۳-۳ k - مرکز با داده‌های پرت

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] M. Sipser. Introduction to the Theory of Computation. Cengage Learning, 2012.
- [5] I. Dinur and S. Safra. On the hardness of approximating minimum vertex cover. Annals of mathematics, pages 439–485, 2005.
- [6] V. V. Vazirani. Approximation algorithms. Springer Science & Business Media, 2013.
- [7] 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.
- [8] R. G. Michael and S. J. David. Computers and intractability: a guide to the theory of np-completeness. WH Freeman & Co., San Francisco, 1979.
- [9] N. Megiddo and K. J. Supowit. On the complexity of some common geometric location problems. SIAM Journal on Computing, 13(1):182–196, 1984.
- [10] M. Bern and D. Eppstein. Approximation algorithms for geometric problems. Approximation algorithms for NP-hard problems, pages 296–345, 1996.

- [11] J. Han, M. Kamber, and J. Pei. Data mining: concepts and techniques: concepts and techniques. Elsevier, 2011.
- [12] P. K. Agarwal and C. M. Procopiuc. *Exact and approximation algorithms for clustering*. Algorithmica, 33(2):201–226, 2002.
- [13] N. Megiddo. *On the complexity of some geometric problems in unbounded dimension*. Journal of Symbolic Computation, 10(3):327–334, 1990.
- [14] B. Chazelle and J. Matoušek. *On linear-time deterministic algorithms for optimization problems in fixed dimension*. Journal of Algorithms, 21(3):579–597, 1996.
- [15] T. M. Chan. *More planar two-center algorithms*. Computational Geometry: Theory and Applications, 13(3):189–198, 1999.
- [16] P. K. Agarwal, R. B. Avraham, and M. Sharir. *The 2-center problem in three dimensions*. Computational Geometry, 46(6):734–746, 2013.
- [17] H. Zarrabi-Zadeh. *Core-preserving algorithms*. In Proceedings of the 20th Canadian Conference on Computational Geometry, pages 159–162, 2008.
- [18] 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.
- [19] 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.
- [20] S. Guha. *Tight results for clustering and summarizing data streams*. In Proceedings of the 12th International Conference on Database Theory, pages 268–275, 2009.
- [21] 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.
- [22] 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.

-
- [23] 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.
- [24] 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.
- [25] S.-S. Kim and H.-K. Ahn. *An improved data stream algorithm for clustering*. In Proceedings of the 11th, pages 273–284. 2014.

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