



هندسه محاسباتی پیشرفته

## مساله نزدیک‌ترین همسایه متنوع

سپیده آقاملائی

۲۰ بهمن ۱۳۹۲

### چکیده

گزارش کردن  $k$  نزدیک‌ترین همسایه متنوع در جستجوی آگاه نسبت به تنوع کاربرد دارد. تعریف مساله به این صورت است که یک نقطه‌ی  $q$  داده شده است، مجموعه  $S$  شامل  $k$  نقطه درون توپ به شعاع  $r$  با بیشترین تنوع را گزارش کنید. تنوع یک مجموعه مثل  $S$  با حداقل فاصله بین دو نقطه‌ی  $S$  سنجیده می‌شود که هرچه بیشتر باشد بهتر است. در این مقاله دو الگوریتم تقریبی برای حالتی که نقاط درون یک فضای فاصله‌ی همینگ  $d$ -بعدی قرار دارند، ارائه می‌شود. زمان کوثری الگوریتم خطی نسبت به  $n$  و چندجمله‌ای بر حسب پارامتر تنوع  $k$  و تعداد ابعاد  $d$  است. برای مقادیر کم  $k$ ، الگوریتم زمان کوثری خطی دارد، حتی اگر تعداد نقاط با فاصله‌ی  $r$  از نقطه کوثری  $q$  بر حسب  $n$  خطی باشند. این الگوریتم‌ها اولین الگوریتم‌هایی از این نوع هستند که پیچیدگی قابل اثبات دارند.

### ۱ مقدمه

مساله‌ی گزارش همسایه‌های نزدیک (یا جستجوی بازه‌ای) به این صورت تعریف می‌شود: مجموعه  $P$  از  $n$  نقطه داده شده است، ساختمان داده‌ای بسازید که به ازای هر نقطه کوثری داده شده، همه‌ی نقاط داده با فاصله کمتر از  $r$  از نقطه کوثری را گزارش کند. این مساله در زمینه‌های مختلفی از جمله پایگاه داده، داده‌کاوی، ذخیره و بازیابی اطلاعات، پایگاه داده‌های تصویر و فیلم، تشخیص الگو، آمار و تحلیل داده‌ها کاربرد دارد. در این کاربردها ویژگی‌های هر شیء مورد بحث (متن، تصویر، ...) به صورت یک نقطه در فضای  $d$  بعدی بازنمایی می‌شوند و ملاک فاصله برای اندازه‌گیری شباهت اشیا به کار می‌رود. در نتیجه مساله اصلی اندیس‌گذاری یا جستجوی معنایی برای اشیا کوثری است. تعداد ویژگی‌ها (ابعاد) از ده تا هزاران متغیر است. یکی از مشکلات اصلی جستجوی معنایی تعداد جواب‌هایی است که بازیابی و گزارش می‌شوند. اگر اندازه

مجموعه جواب خیلی کوچک باشد (مثلاً تنها چند نقطه که به کوثری نزدیک تر اند)، جواب می تواند بسیار همگن و فاقد اطلاعات مفید باشد. اگر تعداد نقاط گزارش شده خیلی زیاد باشد، زمان بازیابی آنها زیاد می شود. همچنین جواب های بزرگ هم خیلی اطلاعات مفیدی نمی دهند. در چند سال اخیر، این مشکل هدف تحقیقات زیادی روی جستجوهای آگاه نسبت به تنوع بوده است. هدف این کار طراحی الگوریتم های کارا برای بازیابی جوابهای مرتبط (نزدیک به نقطه کوثری) و متنوع است. تنوع می تواند به شکل های مختلفی تعریف شود. یکی از روش های مشهور خوشه بندی جوابها و برگرداندن مرکزها است. این روش می تواند زمان اجرای زیاد داشته باشد اگر نقاط تعدادشان زیاد باشد.

## ۱.۱ نتایج ما

در این مقاله دو الگوریتم تقریبی کارا برای مساله  $k$  همسایه نزدیک متنوع داده می شود. مساله به این شکل تعریف می شود: به ازای یک نقطه کوثری داده شده، مجموعه نقاط متنوع  $S$  از  $k$  نقطه را در توپ با شعاع  $r$  حول  $q$  گزارش کنید. تنوع مجموعه  $S$  کمینه ی فاصله ی بین هر دو نقطه ی  $S$  است. به عبارت دیگر، الگوریتم جواب تقریبی مساله خوشه بندی  $k$ -مرکز روی لیست نقاط نزدیک به کوثری را گزارش می کند. زمان اجرا، ضریب تقریب و فضای حافظه ی الگوریتم در جدول داده شده اند. توجه کنید الگوریتم  $A$  از الگوریتم  $B$  بدتر است اما الگوریتم ساده تری است و پیاده سازی و تحلیل آن هم ساده تر است و قبلاً از آن برای بازیابی اخبار متنوع استفاده شده است.

ویژگی های اصلی الگوریتم ما تضمین زمان کوثری برای خطی بودن بر حسب  $n$  و چندجمله ای بودن بر حسب پارامتر تنوع  $k$  و تعداد ابعاد  $d$  است، در حالی که ضریب تقریب ثابت برای رسیدن به تنوع تامین شود. توجه کنید برای مقادیر کم  $k$  الگوریتم ما زمان کوثری خطی دارد، حتی اگر تعداد نقاط به فاصله ی حداکثر  $r$  از  $q$  بر حسب  $n$  خطی باشند. توجه کنید تقریب زدن تنوع اجتناب ناپذیر است چون پیدا کردن مجموعه با اندازه ی  $k$  که تنوع را بیشینه کند با ضریب کمتر از  $2$ ،  $NP$ -سخت است. اینها اولین الگوریتم هایی از این نوع هستند که تضمین قابل اثبات می دهند. یکی از الگوریتم های ما (الگوریتم  $A$ ) به الگوریتم هایی که در عمل برای این کار استفاده می شوند نزدیک است. هرچند این مقالات تضمینی برای کیفیت جواب نداده اند.

## ۲.۱ کارهای گذشته

در این قسمت مروری از کارهای گذشته روی تقریب همسایه نزدیک و جستجوی آگاه نسبت به تنوع که به نتایج این مقاله مرتبط هستند، انجام می دهیم.

همسایه نزدیک. مساله همسایه نزدیک موضوع تحقیقات بسیاری بوده است. چندین الگوریتم کارا برای وقتی که تعداد ابعاد  $d$  کم باشد وجود دارد. با این وجود جوابهای فعلی یا زمان کوثری یا حافظه ی نمایی بر حسب  $d$  دارند. بنابراین در سالهای اخیر روشهای تقریبی برای از بین بردن این تنگنا پیشنهاد شده اند. در تقریب گزارش همسایه نزدیک/جستجوی بازه ای، الگوریتم باید همه ی نقاط در فاصله  $r$  از  $q$  را برگرداند و می تواند تعدادی نقطه در فاصله  $cr$  از  $q$  را هم برگرداند.

یکی از روشهای مشهور برای نزدیک ترین همسایه در ابعاد بالا بر مبنای مفهوم locality sensitive hashing (LSH) است. ایده هش کردن نقاط توسط تعدادی ( $L$ ) تابع هش است تا مطمئن شویم برای هر تابع، احتمال برخورد برای نقاط نزدیک به هم بسیار بیشتر از نقاط دور از هم است. سپس می توان گزارش همسایه نزدیک را با هش کردن نقطه ی کوثری و بازیابی تمام اعضای که در سطل های شامل آن نقطه هستند؛ حل کرد. این روش مثلاً در بسته ی E2LSH برای جستجوی معنایی ابعاد بالا استفاده شده

است. الگوریتم  $LSH$  بر حسب تابع‌های فاصله‌ی به کار رفته در آن انواع مختلفی دارد. در ساده‌ترین حالت وقتی عدم شباهت بین نقاط کوثری با فاصله‌ی همینگ تعریف شود، الگوریتم تضمین ۱) هر نقطه در فاصله کمتر از  $r$  از  $q$  با یک احتمال ثابت (قابل تنظیم) گزارش می‌شود. ۲) زمان کوثری حداکثر  $O(d(n^{1/c} + P_{cr}(q)))$  است که  $P_{cr}(q)$  مجموعه نقاطی از  $P$  را مشخص می‌کند که تا  $q$  فاصله‌ی  $R$  دارند. پس اگر اندازه‌ی مجموعه جواب  $P_{cr}(q)$  بزرگ باشد، کارایی الگوریتم کاهش پیدا می‌کند. به صورت اکتشافی افزایش سرعت می‌تواند با خوشه‌بندی نقاط در هر سطل و نگه داشتن فقط مرکز خوشه‌ها انجام شود. هر چند الگوریتم حاصل هیچ تضمینی نداشته است. (تا کنون) تنوع.

در این مقاله ما تعریف  $content - based$  برای تنوع را به کار بردیم. این روش گزارش کردن  $k$  جواب است که به مقدار کافی از هم متفاوت هستند. این به صورت بیشینه کردن کمترین فاصله بین هر زوج جواب است یا میانگین فاصله بین جواب‌ها و غیره. در این مقاله ما از کمترین فاصله برای بیان ریاضی استفاده می‌کنیم و از الگوریتم خوشه‌بندی حریصانه برای پیدا کردن  $k$  خوشه‌بندی با بیشترین تنوع نقاط در یک مجموعه داده شده استفاده می‌کنیم.

تنها کار قبلی که مستقیماً تعریف ما از  $k$  همسایه نزدیک متنوع را ارائه داده است (۱) است. این مقاله الگوریتمی (مشابه الگوریتم  $A$  در این مقاله ارائه می‌دهد، اما از اندیس جا‌کارد به جای فاصله همینگ استفاده می‌کند) و آن را روی مسائلی در بازیابی اخبار اجرا می‌کند. هرچند این مقاله هیچ تضمینی روی دقت جوابهای گزارش شده نمی‌دهد.

### ۳.۱ تکنیک ما

هر دوی الگوریتم‌هایی که ما استفاده می‌کنیم از  $LSH$  به عنوان پایه استفاده می‌کنند. مشکل اصلی کاهش وابستگی زمان کوثری به اندازه‌ی مجموعه‌ی  $P_{cr}(q)$  از نقاط نزدیک به  $q$  است. الگوریتم اول ( $A$ ) به این صورت این مشکل را حل می‌کند که فقط  $k$  متنوع‌ترین نقاط را در هر سطل نگه می‌دارد. این کار تضمین می‌کند کل تعداد نقاطی که در زمان هر کوثری بررسی می‌شوند حداکثر  $O(kL)$  است که  $L$  تعداد توابع هش است. اما اثبات تضمین‌های تقریب برای این الگوریتم نیازمند این است که هیچ داده‌ی دورافتاده‌ای (نقطه‌ای با فاصله‌ی بیشتر از  $cr$  از  $q$ ) در هیچ سطلی ذخیره نشده باشد. در غیر این صورت این نقطه می‌تواند به عنوان یکی از  $k$  نقطه‌ی متنوع از آن سطل گزارش شود و جایگزین یک نقطه‌ی معتبر شود. این نیاز باعث می‌شود الگوریتم تنها در صورتی کار کند که ضریب تقریب فاصله ( $c$ ) از ۲ بزرگتر باشد.

تضمین ضریب تقریب ۶ برای تنوع با استفاده از مفهوم مجموعه‌های هسته نشان داده شده‌است. به سادگی می‌فهمیم بیشینه‌ی  $k$ -متنوع بودن یک مجموعه نقطه کمتر از ۲ برابر هزینه‌ی جواب بهینه‌ی خوشه‌بندی ( $k - 1$ ) -مرکز آن است. برای مساله‌ی دومی روش ساختن زیرمجموعه‌از نقاط ورودی (یک مجموعه هسته) شناخته شده است که برای هر مجموعه‌ای از مرکزهای خوشه‌ها هزینه‌ی خوشه‌بندی مجموعه هسته، ضریب ثابتی از جواب بهینه‌ی خوشه‌بندی کل داده‌ها باشد. الگوریتم ما سپس به سادگی تنها یک مجموعه‌ی هسته به ازای هر سطل  $LSH$  محاسبه و ذخیره می‌کند. مشخصات استاندارد مجموعه هسته نتیجه می‌دهند که اجتماع مجموعه‌های هسته‌ی سطل‌های متناظر یک نقطه کوثری  $q$  یک مجموعه‌ی هسته برای همه‌ی نقاط این سطل‌ها است. پس اجتماع همه‌ی مجموعه‌های هسته اطلاعات کافی برای بازیابی یک جواب تقریبی بهینه برای همه‌ی نقاط نزدیک به  $q$  می‌دهد.

برای به دست آوردن الگوریتمی که برای  $c > 1$  دلخواه کار کند، باید الگوریتم ما نسبت به داده‌های دور افتاده حساس نباشد ( $robust$  باشد). تحلیل روش  $LSH$  استاندارد تضمین می‌کند که تعداد نقاط دورافتاده در همه‌ی سطل‌ها حداکثر  $O(L)$  است. الگوریتم  $B$  به ویژگی  $robustness$  با نگهداری یک مجموعه هسته‌ی  $robust$  دست پیدا می‌کند که می‌تواند تعدادی از داده‌های دور افتاده را تحمل کند. چون ما از قبل نمی‌دانیم چه تعداد داده‌ی دورافتاده‌ای در یک سطل داده‌شده هست، الگوریتم ما دنباله‌ای از نقاط

را نگه می‌دارد که یک مجموعه هسته‌ی *robust* نسبت به تعداد داده‌های دورافتاده‌ی در حال افزایش را بازنمایی می‌کنند. در زمان کوثری زدن، الگوریتم لیست را پیمایش می‌کند تا نقاط کافی برای اطمینان از *robustness* خوانده شده باشند.

## ۲ تعریف مساله

فرض کنید  $(\Delta, dist)$  یک فضای متریک  $d$  بعدی باشد. با دو تعریف زیر شروع می‌کنیم.

**تعریف ۱** به ازای هر مجموعه  $S \in \Delta$  داده شده، تنوع به صورت کمینه فاصله‌ی بین هر دو نقطه از نقاط مجموعه تعریف می‌شود:

$$div(S) = \min_{p, p' \in S} dist(p, p')$$

**تعریف ۲** به ازای هر مجموعه  $S \in \Delta$  داده شده،  $k$ -تنوع آن بیشینه‌ی تنوع قابل دستیابی با انتخاب یک زیرمجموعه  $k$  عضوی است:

$$div_k(S) = \max_{S' \subset S, |S'|=k} div(S')$$

ما همچنین به مجموعه‌ی  $S'$  که مقدار بیشینه‌ی رابطه‌ی بالا را می‌دهد  $k$ -زیرمجموعه بهینه‌ی  $S$  می‌گوییم. توجه کنید  $k$ -تنوع در حالتی که  $|S| < k$  باشد تعریف نشده است.

برای اجتناب از  $k$ -تنوع مجموعه‌های با اعضای کمتر از  $k$ ، قرارداد می‌کنیم هر نقطه‌ی  $p$  از مجموعه نقاط ورودی  $P$  را  $k$  بار تکرار کنیم. این کار باعث می‌شود همهی مجموعه‌های ناتهی  $S$  که در ادامه‌ی مقاله آمده‌اند، کیفیت  $div_k(S)$  خوش‌تعریف و اگر تعداد نقاط کمتر از  $k$  باشد،  $\bullet$  داشته باشند. به سادگی می‌فهمیم که این کار کرانه‌های حافظه‌ی الگوریتم را تغییر نمی‌دهد.

مساله‌ی همسایه نزدیک  $k$ -متنوع به این صورت تعریف می‌شود: به ازای نقطه‌ی کوثری داده شده‌ی  $q$ ، مجموعه‌ی  $S$  را طوری گزارش کنید که: (۱)  $S \in P \cap B(q, r)$  باشد که  $B(q, r) = \{p \mid dist(p, q) \leq r\}$  توپ با شعاع  $r$  و مرکز  $q$  است. (۲)  $|S| = k$  (۳)  $div(S)$  بیشینه باشد.

به دلیل اینکه الگوریتم ما تقریبی است، باید مساله‌ی همسایه نزدیک  $k$ -متنوع تقریبی را تعریف کنیم. در این حالت باید برای یک ضریب تقریب  $c > 1$  و  $\alpha > 1$  داشته باشیم: (۱)  $S \subset P \cap B(q, cr)$  (۲)  $div(S) \geq \frac{1}{\alpha} div_k(P \cap B(q, r))$  (۳)  $|S| = k$ .

## ۳ پیشنهادها

### ۱.۳ الگوریتم GMM

فرض کنید مجموعه نقاط  $S \subset \Delta$  را داریم و می‌خواهیم یک  $k$ -زیرمجموعه بهینه برای  $S$  محاسبه کنیم. یعنی مجموعه‌ای از  $k$  نقطه پیدا کنیم که فاصله‌ی دو به دوی آنها بیشینه شود. هرچند این مساله  $np$ -سخت است، یک الگوریتم ساده‌ی ۲-تقریبی حریصانه به نام *GMM* برای آن وجود دارد. در این مقاله ما از یک نسخه‌ی کمی تغییر یافته‌ی *GMM* استفاده می‌کنیم (اثبات ضریب تقریب این نسخه دقیقاً مثل همان اثبات اصلی است). به الگوریتم یک مجموعه نقاط  $S$  و یک پارامتر  $k$  به عنوان ورودی داده می‌شود. ابتدا الگوریتم یک نقطه‌ی دلخواه  $a \in S$  انتخاب می‌کند. سپس به صورت تکراری نقطه‌ی بعدی را به مجموعه خروجی اضافه می‌کند تا  $k$  نقطه داشته باشد. به طور دقیقتر، در هر قدم، به صورت

حریصانه نقطه‌ای که کوتاهترین فاصله‌اش تا نقاط انتخاب شده فعلی بیشینه باشد، اضافه می‌کند. توجه کنید که تغییر اینکه همه‌ی نقاط  $k$  بار تکرار شوند باعث می‌شود که اگر مجموعه‌ی ورودی کمتر از  $k$  نقطه متمایز داشته باشد، مجموعه‌ی خروجی  $S'$  شامل همه‌ی نقاط باشد.

### الگوریتم ۱ GMM

- ۱: کامنت کار نمی‌کرد:
- ۲: ورودی  $S$ : مجموعه نقاط،  $k$ : اندازه‌ی زیرمجموعه
- ۳: خروجی  $S'$ : زیرمجموعه‌ای از  $S$  با اندازه‌ی  $k$
- ۴: یک نقطه‌ی دلخواه:  $a \leftarrow S'$
- ۵: به ازای  $i$  از ۲ تا  $k$ :
- ۶:  $p \in S - S'$  پیدا کن که  $\min_{x \in S'} \text{dist}(p, x)$  را بیشینه کند.
- ۷:  $S' \leftarrow S' \cup \{p\}$
- ۸:  $S'$  را برگردان.

لم ۱ زمان اجرای الگوریتم  $O(k \cdot |S|)$  است و ضریب تقریب حداکثر ۲ برای  $k$ -تنوع  $\text{div}_k(S)$  دارد.

### ۲.۳ مجموعه‌های هسته

تعریف ۳ فرض کنید  $(P, \text{dist})$  یک متریک باشد. برای هر زیرمجموعه‌ای از نقاط  $S, S' \subset P$ ، هزینه‌ی  $k$ -مرکز را به صورت زیر تعریف می‌کنیم:

$$KC(S, S') = \max_{p \in S} \min_{p' \in S'} \text{dist}(p, p')$$

مساله‌ی  $k$ -مرکز متریک به این صورت تعریف می‌شود:  $S$  داده شده‌است، زیرمجموعه‌ی  $S' \in S$  با اندازه‌ی  $k$  پیدا کنید که  $KC(S, S')$  را کمینه کند. ما این هزینه‌ی بهینه را با  $KC_k(S)$  نشان می‌دهیم.

لم ۲  $k$ -تنوع یک مجموعه  $S$  ارتباط نزدیکی با هزینه‌ی بهترین  $(k-1)$ -مرکز  $S$  دارد. یعنی:

$$KC_{k-1}(S) \leq \text{div}_k(S) \leq 2KC_{k-1}(S)$$

اثبات. برای نامساوی اول فرض کنید  $S'$   $k$ -زیرمجموعه بهینه‌ی  $S$  باشد. همچنین  $a \in S'$  را یک نقطه‌ی دلخواه و  $S' - a = S' - \{a\}$ . سپس برای هر نقطه‌ی  $b \in S - S'$ ، داریم  $\min_{p \in S'} \text{dist}(b, p) \leq \text{dist}(b, a)$ ، یعنی  $\text{div}(b \cup S' - a) > \text{div}(S')$  پس  $KC(S, S' - a) \leq \text{div}_k(S)$  است و نامساوی برقرار است. برای قسمت دوم، فرض کنید  $C = \{a_1, \dots, a_{k-1}\}$  یک مجموعه بهینه از  $(k-1)$ -مرکز برای  $S$  باشد. سپس چون  $S'$  اندازه‌ی  $k$  دارد، طبق اصل لانه‌کبوتری،  $p, p' \in S'$  و  $a$  وجود دارند که

$$a = \arg \min_{c \in C} \text{dist}(p, c) = \arg \min_{c \in C} \text{dist}(p', c)$$

و در نتیجه طبق نامساوی مثلث داریم:

$$\text{div}_k(S) = \text{div}(S') \leq \text{dist}(p, p') \leq \text{dist}(p, a) + \text{dist}(a, p') \leq 2KC_{k-1}(S)$$

□

**تعریف ۴** فرض کنید  $(P, dist)$  متریک ما باشد. پس برای  $\beta \geq 1$  تعریف می‌کنیم: یک مجموعه‌ی هسته‌ی  $\beta$ ، برای مجموعه نقاط  $S \subset P$  هر زیرمجموعه‌ی  $S' \subset S$  باشد که برای هر زیرمجموعه‌ی  $(k-1)$  نقطه‌ی  $F \subset P$  داشته باشیم:  $KC(S', F) \geq \beta KC(S, F)$ .

**تعریف ۵** فرض کنید  $(P, dist)$  متریک ما باشد. پس برای هر  $\beta \leq 1$  و عدد صحیح  $l$ ، یک  $l - robust\beta - coresets$  برای هر مجموعه نقطه‌ی  $S \subset P$  را به این صورت تعریف می‌کنیم: زیرمجموعه‌ی  $S' \subset S$  طوری که به ازای هر مجموعه نقاط دورافتاده‌ی  $O \subset P$  با حداکثر  $l$  نقطه،  $S' - O$  یک  $\beta -$  مجموعه هسته‌ی  $S - O$  باشد.

### ۳.۳ هشینگ حساس به محل

هشینگ حساس به محل روشی برای حل مسائل همسایه نزدیک تقریبی است. ایده‌ی اصلی آن هش کردن داده و نقطه کوثری به صورتی است که احتمال برخورد برای نقاطی که نزدیک به هم هستند از آنهایی که از هم دور هستند بسیار بیشتر باشد. تعریف ریاضی آن در ادامه آمده است.

**تعریف ۶** یک خانواده‌ی  $H = h : \Delta \rightarrow U$  را برای  $(\Delta, dist)$ ،  $(r_1, r_2, p_1, p_2)$ -حساس می‌گویند، اگر برای هر  $p, q \in \Delta$  داشته باشیم:

• اگر  $dist(p, q) \leq r_1$  باشد، آنگاه  $Pr_H[h(q) = h(p)] \geq p_1$  باشد.

• اگر  $dist(p, q) \leq r_2$  باشد، آنگاه  $Pr_H[h(q) = h(p)] \leq p_2$  باشد.

برای اینکه یک خانواده‌ی حساس به محل مفید باشد، باید در نامساوی‌های  $p_1 > p_2$  و  $r_1 < r_2$  صدق کند.

برای یک خانواده‌ی  $LSH$  داده شده، الگوریتم  $L$  تابع هش  $g_1, \dots, g_L$  و آرایه‌های متناظر آنها یعنی  $A_1, \dots, A_L$  را تولید می‌کند. هر تابع هش به شکل  $h_i = \langle h_{i,1}, \dots, h_{i,K} \rangle$  است که  $h_{i,j}$  به صورت تصادفی یکنواخت از  $H$  انتخاب می‌شود. سپس هر نقطه‌ی  $p$  در یک سطل  $g_i(p)$  از  $A_i$ ها ذخیره می‌شود ( $1 \leq i \leq L$ ). برای پاسخگویی به یک کوثری  $q$ ، ما در  $(A_L(g_L(q)) \cup \dots \cup A_1(g_1(q)))$  جستجو می‌کنیم. یعنی برای هر آرایه تنها یک سطل که متناظر نقطه کوثری  $q$  است را بررسی می‌کنیم. در این مقاله برای سادگی  $LSH$  را برای فاصله‌ی همینگ در نظر می‌گیریم. هرچند نتایج مشابه برای توابع  $LSH$  عمومی قابل اثبات هستند. لم زیر را از (۱۰) مرور می‌کنیم.

**لم ۳** فرض کنید  $dist(p, q)$  همینگ باشد برای  $p, q \in \Sigma^d$  که  $\Sigma$  هر الفبای متناهی باشد. پس برای هر  $r, c \geq 1$  وجود دارد خانواده‌ی  $H$  که  $(r, rc, p_1, p_2)$ -حساس هستند که  $p_1 = 1 - r/d$  و  $p_2 = 1 - rc/d$ . همچنین اگر قرار دهیم  $\rho = \frac{\log 1/p_1}{\log 1/p_2}$  آنگاه داریم  $\rho \leq 1/c$  است. به علاوه با اضافه کردن صفرهای قبل از عدد می‌توانیم فرض کنیم  $r/d \leq 1/2$ .

## ۴ الگوریتم $A$

الگوریتم  $A$  که اولین بار در (۱) ارائه شد، بر مبنای الگوریتم  $LSH$  است. در زمان پیش‌پردازش،  $LSH$  تابع هش  $g_1, \dots, g_L$  و آرایه‌های  $A_1, \dots, A_L$  را می‌سازد. بعد هر نقطه  $p$  در سطل‌های  $A_i[g_i(p)]$  ذخیره می‌شود (به ازای همه‌ی  $1..L$ ). به علاوه برای هر آرایه  $A_i$  الگوریتم از  $GMM$  استفاده می‌کند

تا ۲- تقریبی از  $k$ - زیرمجموعه بهینه برای هر سطل پیدا کند و آن را در سطل متناظر  $A'_i$  ذخیره می‌کند. این زیرمجموعه یک  $1/3$  هسته برای نقاط سطل است. یک نقطه کوثری  $q$  داده شده است، الگوریتم اجتماع سطل‌ها  $Q = A'_1(g_1(q)) \cup \dots \cup A'_L(g_L(q))$  را حساب می‌کند و بعد از آن همه‌ی نقاط دورافتاده (یعنی نقاطی که فاصله‌شان تا  $q$  بیشتر از  $cr$  است) را حذف می‌کند. در گام آخر الگوریتم  $GMM$  را روی  $Q$  اجرا می‌کند و جواب تقریباً بهینه‌ی  $k$ - زیرمجموعه‌ی  $Q$  را برمی‌گرداند. شبه کدها در ادامه آمده‌اند و در بخش بعد به تحلیل درستی الگوریتم می‌پردازیم.

### الگوریتم ۲ پیش پردازش

- ۱: ورودی:  $G = \{g_1, \dots, g_L\}$ : مجموعه‌ی  $L$  تابع هش،  $P$ : مجموعه نقاط،  $k$
- ۲: خروجی:  $A' = \{A'_1, \dots, A'_L\}$
- ۳: به ازای همه نقاط  $p \in P$ :
- ۴: به ازای همه‌ی توابع هش  $g_i \in G$ :
- ۵: نقطه‌ی  $p$  را به سطل  $A_i[g_i(p)]$  اضافه کن.
- ۶: به ازای  $A_i \in A$ :
- ۷: به ازای  $j = 1 \rightarrow size(A_i)$ :
- ۸: کامنت: فقط نقاط تقریب  $k$ -متنوع رادر هر سطل نگه دار.
- ۹:  $A'_i[j] = GMM(A_i[j], k)$

### الگوریتم ۳ پردازش کوثری

- ۱: کامنت: ورودی  $q$ : نقطه کوثری،  $k$
- ۲: خروجی  $Q$ : مجموعه نقاط  $k$ -متنوع
- ۳:  $Q \leftarrow \emptyset$
- ۴: به ازای  $i = 1 \rightarrow L$ :
- ۵:  $Q \leftarrow Q \cup A'_i[g_i(q)]$
- ۶: به ازای همه‌ی  $q \in Q$ :
- ۷: اگر  $dist(p, q) > cr$ :
- ۸:  $p$  را از  $Q$  حذف کن. کامنت: چون این نقطه دورافتاده است.
- ۹:  $Q \leftarrow GMM(Q, k)$
- ۱۰:  $Q$  را برگردان.

## ۱.۴ تحلیل

در این قسمت، ابتدا مقادیر پارامترهای  $L$  و  $K$  را بر حسب  $n$  و  $1/c \leq \rho$  به دست می‌آوریم، طوری که با احتمال ثابت الگوریتم کار کند. اینجا  $L$  تعداد کل توابع هش به کار رفته و  $K$  تعداد توابع هش  $h_{i,j}$  است که در  $g_i$  استفاده شده‌است. همچنین ما باید نشان دهیم محدود کردن اندازه‌ی سطل‌ها به  $k$  متنوع‌ترین نقطه در  $A'$ ، برای به دست آوردن تقریب خوب کارآمد است. این موضوع را در ادامه بحث می‌کنیم.

لم ۴ به ازای  $c > 2$ ، توابع هش  $g_1, \dots, g_L$  به شکل  $g_i = \langle h_{i,1}, \dots, h_{i,K} \rangle$  وجود دارند که  $h_{i,j} \in H$  است که  $H$  و  $p_1$  و  $p_2$  در لم ۳ تعریف شده‌اند، طوری که با قرار دادن  $L = (\log(\mathfrak{K})/p_1)^{1/(1-p)} \times$  دو پیشامد زیر با احتمال ثابت رخ دهند:

•  $\forall p \in Q^* : \exists i, s.t., p \in A_i[g_i(q)]$  که جواب بهینه  $Q^*$  (ک-زیرمجموعه بهینه برای  $P \cap B(q, r)$  است.

•  $\forall p \in \cup_i A_i[g_i(q)] : dist(p, q) \leq cr$  یعنی هیچ داده‌ی دورافتاده‌ای به ازای تمام توابع هش، بین نقاطی که به سطل شامل  $q$  هش می‌شوند، نباشد.

اثبات آن در پیوست A آمده است.

قضیه‌ی ۱ چون هر نقطه یک بار در هر تابع هش، هش می‌شود؛ کل فضای که توسط این الگوریتم استفاده می‌شود حداکثر  $n$  است:

$$nL = n(\log(\mathfrak{K})/p_1(\mathfrak{K}n)^\rho)^{1/(1-p)} = O\left(\left(\frac{n \log k}{1-r/d}\right)^{\frac{1}{1-p}}\right) = O((n \log k)^{1+\frac{1}{c-1}})$$

که برای به دست آوردن آن از این حقیقت که  $1/2 \leq r/d \leq 1/c$ ،  $c > 2$  است، استفاده کرده‌ایم. همچنین به  $O(nd)$  فضا برای ذخیره‌ی نقاط نیاز داریم.

قضیه‌ی ۲ زمان کوثری  $O((\log n)/r + k^2) \cdot (\log k)^{\frac{c}{c-1}} \cdot n^{\frac{1}{c-1}} d$  است.

اثبات. زمان کوثری الگوریتم برای هر کوثری حداکثر  $O(L)$  محاسبه‌ی هش است که هر کدام  $O(K)$  زمان می‌گیرد:

$$\begin{aligned} O(KL) &= O((\log_{1/p_1}(\mathfrak{K}nL)).L) = O\left(\frac{\log n}{\log(1/p_1)} L\right) \\ &= O\left(\frac{d}{r} \log n \cdot \left(\frac{\log k}{1-r/d}\right)^{\frac{c}{c-1}} \cdot n^{\frac{1}{c-1}}\right) \\ &= O\left(\frac{d}{r} (\log k)^{\frac{c}{c-1}} n^{\frac{1}{c-1}} \log n\right) \end{aligned}$$

که برای به دست آوردن آن از تقریب  $1/2 \leq r/d \leq 1/c$ ،  $c \geq 2$ ،  $\log p_1 \approx 1 - p_1 = \frac{cr}{d}$  استفاده کرده‌ایم. همچنین در قدم آخر، لازم است الگوریتم  $GMM$  را حداقل برای متوسط  $kL$  نقطه اجرا کنیم که زمان زیر را لازم دارد:

$$\begin{aligned} O(k^2 Ld) &= O(k^2 \cdot (\log k/p_1(\mathfrak{K}n)^\rho)^{\frac{1}{1-p}} d) \\ &= O(k^2 (\log k)^{\frac{c}{c-1}} \cdot n^{\frac{1}{c-1}} d) \end{aligned}$$

□

لم ۵  $GMM(S, k)$  یک  $1/3$ -هسته برای  $S$  محاسبه می‌کند.  
expectation<sup>۱</sup>



اثبات. فرض کنید مجموعه  $k$  نقطه محاسبه شده توسط  $GMM$  را  $S'$  بنامیم. حالا هر زیرمجموعه‌ی  $(k-1)$  عضوی از نقاط  $F \subset P$  را بردارید. طبق اصل لانه کبوتری  $a, b \in S'$  وجود دارند که نزدیک‌ترین نقطه به آنها در  $F$  یکسان باشد؛ یعنی  $c \in F$  وجود داشته باشد که:

$$c = \operatorname{argmin}_{f \in F} \operatorname{dist}(a, f) = \operatorname{argmin}_{f \in F} \operatorname{dist}(b, f)$$

در نتیجه طبق نامساوی مثلث داریم:

$$\operatorname{div}(S') \leq \operatorname{dist}(a, b) \leq \operatorname{dist}(a, c) + \operatorname{dist}(b, c) \leq 2KC(S', F)$$

حالا هر نقطه‌ی  $s \in S$  را بردارید و  $s'$  را نزدیک‌ترین نقطه از اعضای  $S'$  به  $s$  فرض کنید و  $f$  را نزدیک‌ترین نقطه از  $F$  به  $s'$ . همچنین فرض کنید  $a \in S'$  نقطه‌ای باشد که در قدم آخر الگوریتم  $GMM$  اضافه شده است. پس طبق تعریف  $s'$  و  $a$  و انتخاب حریصانه‌ی  $GMM$  داریم:

$$\operatorname{dist}(s, s') \leq \min_{p \in S' \setminus \{a\}} \operatorname{dist}(p, s) \leq \min_{p \in S' \setminus \{a\}} \operatorname{dist}(p, a) \leq \operatorname{div}(S')$$

بنابراین طبق نامساوی مثلث:

$$\operatorname{dist}(s, f) \leq \operatorname{dist}(s, s') + \operatorname{dist}(s', f) \leq \operatorname{div}(S') + KC(S', F) \leq 3KC(S', F)$$

از آنجایی که این برای هر  $s \in S$  برقرار است، می‌توانیم نتیجه بگیریم  $KC(S, F) \leq 3KC(S', F)$  □  
که از اینجا حکم نتیجه می‌شود.

لم ۶ فرض کنید  $S_1, \dots, S_m$  زیرمجموعه‌هایی از  $P$  باشند و  $S = \cup_i S_i$ . همچنین فرض کنید  $T_i = GMM(S_i, k)$  یک  $2$ -تقریب برای  $k$ -زیرمجموعه بهینه برای  $S_i$  باشد که با اجرای الگوریتم  $GMM$  روی  $S_i$  به دست آمده است. همچنین تعریف کنید  $T = \cup_i T_i$  و فرض کنید  $T' = GMM(T, k)$  یک  $2$ -تقریب برای  $k$ -زیرمجموعه بهینه‌ی  $T$  باشد. در این صورت داریم:  $\operatorname{div}(T') \geq \frac{1}{6} \operatorname{div}_k(S)$

اثبات. فرض کنید  $S'$   $k$ -زیرمجموعه بهینه‌ی  $S$  باشد. همچنین فرض کنید  $a \in T'$  آخرین نقطه‌ای باشد که توسط الگوریتم  $GMM(T, k)$  اضافه شده است و قرار دهید  $a \in T' -$ . طبق اصل لانه کبوتری  $c \in T' -$  و  $p, p' \in S'$  وجود دارند که:

$$c = \operatorname{argmin}_{t \in T' -} \operatorname{dist}(t, p) = \operatorname{argmin}_{t \in T' -} \operatorname{dist}(t, p')$$

پس طبق نامساوی مثلث، لم ۵ و این حقیقت که اجتماع  $\beta$ -هسته‌های  $S_i$  یک  $\beta$ -هسته برای  $S$  است، داریم:

$$\begin{aligned} \operatorname{div}_k(S) &= \operatorname{div}(S') \leq \operatorname{dist}(p, p') \leq \operatorname{dist}(p, c) + \operatorname{dist}(p', c) \\ &\leq 2KC(S, T' -) \leq 6KC(T, T' -) \end{aligned}$$

و چون برای هر  $b \in T$  داریم:

$$\min_{t \in T' -} \operatorname{dist}(t, b) \leq \min_{t \in T' -} \operatorname{dist}(t, a) \leq \operatorname{div}(T')$$

در نتیجه  $KC(T, T' -) \leq \operatorname{div}(T')$  و لم اثبات می‌شود. □

قضیه‌ی ۳ نتیجه: با احتمال ثابت، ضریب تقریب به دست آمده توسط این الگوریتم ۶ است.

اثبات. فرض کنید  $S = \cup_i S_i$  و  $S_i = A_i[g_i(q)]$  باشد. همچنین فرض کنید  $T = T_i = A'_i[g_i(q)]$  باشد. به علاوه  $Q^*$  را  $k$ -زیرمجموعه بهینه  $P \cap B(q, r)$  تعریف کنید و  $T' = GMM(T, k)$  و  $Q$  را مجموعه خروجی فرض کنید. از توضیح الگوریتم، بدیهی است که  $Q \subset B(q, cr)$  است. پس تنها چیزی که باید ثابت کنیم تنوع است. طبق قضیه ۴، با احتمال ثابت دو گزاره‌ی زیر برقرار هستند:

•  $S \subset B(q, cr)$  بنابراین داریم  $T \subset B(q, cr)$  که نشان می‌دهد وقتی الگوریتم ۳ را روی  $q$  اجرا می‌کنیم، چون  $T$  شامل هیچ داده دورافتاده‌ای نیست، الگوریتم  $GMM$  در خط ۱۰ روی مجموعه  $T$  فراخوانی می‌شود و بنابراین  $Q = T'$ .

•  $Q^* \subset S$  پس داریم  $div_k(Q^*) = div(Q^*)$

پس طبق لم ۶ داریم:

$$div(Q) \geq \frac{1}{c} div_k(S) \geq \frac{1}{c} div(Q^*)$$

□

## ۵ الگوریتم $B$

در این قسمت، نسخه‌ی تغییر یافته‌ی الگوریتم  $A$  را معرفی و تحلیل می‌کنیم که ضریب تقریب ثابت دارد. فرض کنید می‌دانیم تعداد نقاط دورافتاده در هر سطل حداکثر  $l$  است. پس می‌توانیم برای هر سطل آرایه‌ی  $A$ ، یک  $l - robust \ 1/3 - coresets$  در سطل مربوطه در آرایه  $A'$  ذخیره کنیم. ابتدا در الگوریتم ۴ نشان می‌دهیم چطور یک  $l - robust \ \beta - coresets$  را پیدا کنیم اگر بدانیم چطور یک  $\beta - coresets$  را پیدا کنیم. این الگوریتم (۳) است که مجموعه‌های هسته را پوست می‌کند تا  $\beta - coresets$  به دست بیاورد.

**الگوریتم ۴** هسته و داده دورافتاده

۱: کامنت: ورودی  $S$ : مجموعه نقاط

۲: خروجی  $S'$ : آرایه‌ای که یک  $(l, \beta) - coresets$  برای  $S$  است.

۳:  $S' \leftarrow \emptyset$

۴: به ازای  $i = 1 \rightarrow (l + 1)$ :

۵:  $R_i \leftarrow \beta - coresets(S)$

۶:  $R_i$  را به انتهای  $S'$  اضافه کن.

۷:  $S \leftarrow S \setminus R_i$

۸:  $S'$  را برگردان.

لم ۷ فرض کنید  $O \subset P$  یک زیرمجموعه از داده‌های دورافتاده باشد و  $S'_j$  مجموعه‌ی  $(k, j)$  نقطه‌ی  $S'$  باشد که بعد از  $j$ -امین مرحله الگوریتم به دست آمده باشد. در این صورت برای هر  $l \geq j \geq 0$  که در شرط  $|S'_{j+1} \cap O| \leq j$  صدق کند،  $S'_{j+1} \setminus O$  یک  $\beta - هسته$  برای  $S \setminus O$  است.

اثبات. فرض کنید  $F \subset P$  هر زیرمجموعه‌ای از  $(k-1)$  نقطه باشد و  $q$  دورترین نقطه از  $F$  در بین اعضای  $(S \setminus O)$  باشد، یعنی:

$$q = \operatorname{argmax}_{p \in S \setminus O} \min_{f \in F} \operatorname{dist}(p, F)$$

حالا به ازای هر  $i \leq (j+1)$ ، اگر  $q$  نقطه‌ای در  $R_i$  باشد، آنگاه لم برقرار است چون  $KC(S, F) = KC(S'_{j+1}, F)$  در غیر این صورت چون  $q$  در هیچ یک از  $(j+1)$  گام اول انتخاب نشده است، هر کدام از  $R_i$  ها برای  $(i \leq (j+1))$  شامل  $r_i$  هستند که  $KC(\{r_i\}, F) \geq \beta KC(\{q\}, F)$  از بین این  $j+1$  حدافل یکی از آنها در  $O$  نیست، پس  $KC(S'_{j+1} \setminus O, F) \geq \beta KC(S \setminus O, F)$  . □

**قضیه ۴** نتیجه: الگوریتم ۴ به درستی یک  $(l, \beta)$  هسته برای  $S$  محاسبه می‌کند.

اثبات. فرض کنید برای هر مجموعه‌ی  $O \subset P$  از نقاط دورافتاده که  $|O| \leq l$  باشد، شرایط لم ۷ برای  $l = j$  برقرار است. پس وقتی الگوریتم جواب را برمی‌گرداند، یک  $l - \text{robust} \beta - \text{coreset}$  را به درستی محاسبه کرده‌است. □

چون طبق لم ۵، می‌دانیم که  $GMM(S, k)$  یک  $1/3$ -هسته با اندازه‌ی  $k$  برای  $S$  پیدا می‌کند، کافی است خط ۳ الگوریتم ۴ را با  $R \leftarrow GMM(S, k)$  جایگزین کنیم تا به یک  $l - \text{robust} 1/3 - \text{coreset}$  با اندازه‌ی  $k(l+1)$  برای مجموعه  $S$  برسیم.

پس تنها تغییری که در قسمت پیش‌پردازش الگوریتم  $A$  انجام می‌دهیم این است که هر سطل از  $A'_i$  یک  $l - \text{robust} 1/3 - \text{coreset}$  از سطل متناظر  $A_i$  را نگه می‌دارد. پس خط ۸ الگوریتم ۲ به این صورت تغییر می‌کند:  $A'_i[j] = (l, \beta) - \text{coreset}(A_i[j], l = 3L, \beta = 1/3)$ . شبه کد پیش‌پردازش کوثری در الگوریتم ۵ نشان داده شده است. برای هر سطل که متناظر  $q$ ، الگوریتم سعی می‌کند کوچکترین مقدار  $l$  را پیدا کند که تعداد داده‌های دورافتاده در  $(l+1)$  عضو اول  $A'_i[g_i(q)]$  از  $l$  بیشتر نشود. سپس این نقاط را به  $T$  اضافه می‌کند و  $k$  زیرمجموعه بهینه تقریبی از نقاط غیر دورافتاده‌ی  $T$  را برمی‌گرداند.

توجه کنید که حلقه داخلی (خط ۴ تا ۱۱) الگوریتم ۵ می‌تواند به صورت کارا پیاده‌سازی شود. با دانستن تعداد نقاط دورافتاده‌ی  $U_i^j$ ، تنها  $k$  عضو دیگر هستند که باید چک شوند که در  $U_i^{j+1}$  دورافتاده هستند یا نه. همچنین هر نقطه می‌تواند در زمان  $O(d)$  چک شود که دورافتاده هست یا نه. پس کل حلقه‌ی داخلی  $O(kld)$  زمان می‌خواهد.

## ۱.۵ تحلیل

ابتدا قضیه زیر را که شبیه قضیه ۴ است بیان می‌کنیم. اثبات آن خیلی شبیه اثبات الگوریتم اصلی  $LSH$  است که در (۱۰) داده شده است و به همین دلیل حذف شده است.

**قضیه ۵** توابع هش  $g_1, \dots, g_L$  به شکل  $g_i = \langle h_{i,1}, \dots, h_{i,k} \rangle$  وجود دارند که  $h_{i,j} \in H$  که برای  $H$  مقادیر  $p_1$  و  $p_2$  طبق قضیه ۳ با قرار دادن  $n^p/p_1 \times \log(4k)$  و  $L = \lceil \log_{1/p_1} n \rceil$  به دست می‌آیند. با احتمال ثابت دو رویداد زیر رخ می‌دهند:

- $\exists i, s.t., p \in A_i[g_i(q)] \cap P$  که  $Q^*$  جواب بهینه است  $(k - \text{زیرمجموعه بهینه‌ی } P \cap B(q, r))$

## الگوریتم ۵ پردازش کوثری

- ۱: کامنت: ورودی:  $q$  نقطه ورودی
- ۲: خروجی:  $Q$ : مجموعه‌ی نقاط  $k$ -متنوع
- ۳:  $T \leftarrow \emptyset$
- ۴: مجموعه نقاط دورافتاده  $O \leftarrow$
- ۵: به ازای  $i = 1 \leftarrow L$ :
- ۶: به ازای  $l = 0 \leftarrow 3L$ :
- ۷:  $U_i^{l+1} = A_i^l[g_i(q)]$  نقطه‌ی اول از  $k(l+1)$
- ۸: اگر  $|U_i^{l+1} \cap O| \leq l$ :
- ۹:  $l_i \leftarrow l$
- ۱۰:  $T_i \leftarrow U_i^{l+1}$
- ۱۱: از شرط خارج شو.
- ۱۲:  $T \leftarrow T \cup T_i$
- ۱۳:  $Q \leftarrow GMM(T \setminus O, k)$
- ۱۴:  $Q$  را برگردان.

•  $\{p \in \cup_i A_i[g_i(q)] : dist(p, q) > cr\} |le 3L$  یعنی تعداد نقاط دورافتاده از بین نقاط هش شده به سطلی که  $q$  به آن هش شده، حداکثر  $3L$  است.

قضیه‌ی ۶ نتیجه: چون هر نقطه‌ای در هر تابع هش یک بار هش می‌شود و هر سطل  $A_i^l$  زیرمجموعه‌ی متناظر سطل  $A_i$  است، کل فضای مورد نیاز برای الگوریتم حداکثر انقدر است:

$$nL = n \log(3k) \cdot n^p / p_1 = O(\log k \cdot n^{1+1/c})$$

دلیل درستی آن این است که  $ple 1/c$  است و  $1/2 \leq 1 - r/d \leq p_1$  است. همچنین ما  $O(nd)$  فضا برای نگهداری نقاط نیاز داریم.

قضیه‌ی ۷ با احتمال ثابت ضریب تقریب الگوریتم ۶ است.

اثبات. ما ابتدا با مجموعه‌ای از نمادها که برای اثبات لازم است شروع می‌کنیم:

$$S_i = A[g_i(q)], S = \cup_i S_i$$

• فرض کنید  $O$  مجموعه نقاط دورافتاده باشد، یعنی  $O = S \setminus B(q, cr)$ . می‌دانیم با احتمال ثابت  $|O| \leq 3L$

• فرض کنید  $S'$   $k$ -زیرمجموعه بهینه برای  $S \setminus O$  باشد.

•  $U_i = A_i^l[g_i(q)]$  و  $U_i^j$  مجموعه‌ی  $kz$  عضو اول  $U_i$  باشد. (توجه کنید چون خروجی الگوریتم ۴ آرایه است، اعضای آن ترتیب دارند.)  $T_i = U_i^{(l_i+1)}$  تعریف می‌کنیم که  $l_i$  طوری انتخاب شده است که به ازای هر  $l_i < l'_i$  است و داریم  $|U_i^{(l'_i+1)} \cap O| > l'_i$ . این دقیقاً متغیر  $T_i$  در الگوریتم ۵ است. به علاوه  $T = \cup_i T_i$  قرار دهید.

•  $Q^*$  را  $k$ -زیرمجموعه‌ی بهینه‌ی  $P \cap B(q, r)$  تعریف کنید و  $Q$  را خروجی الگوریتم یعنی  $Q = GMM(T \setminus O, k)$  فرض کنید  $a \in Q$  نقطه‌ی اضافه شده در مرحله‌ی آخر الگوریتم  $GMM$  باشد.  $Q - = Q \setminus \{a\}$  تعریف می‌کنیم.

از توضیح الگوریتم بدیهی است  $Q \subset B(q, cr)$ . همچنین طبق قضیه ۱، با احتمال ثابت داریم  $Q^* \subset S$  و تعداد نقاط دورافتاده از  $3L$  بیشتر نیست. پس داریم:  $div_k(Q) = div(Q^*) \leq div_k(S \setminus O)$ . با این شرایط کافی است ثابت کنیم:  $div(Q) \geq div_k(S \setminus O)/6 = div(S')/6$ . طبق اصل لانه کبوتری، چون  $|S'| = k$  و  $|Q - | = k - 1$  پس نقاط  $p, p' \in S'$  وجود دارند که نزدیک‌ترین نقطه به آنها در  $Q -$  یکسان است. یعنی  $c \in Q -$  وجود دارد که  $KC(\{p\}, Q -) = dist(p, c)$  و  $KC(\{p'\}, Q -) = dist(p', c)$  پس طبق نامساوی مثلث داریم:

$$div(S') \leq dist(p, p') \leq dist(p, c) + dist(p', c) \leq 2KC(S \setminus O, Q -)$$

طبق لم ۷،  $T_i \setminus O$  یک  $1/3$ -هسته برای  $S_i \setminus O$  است، پس اجتماع آنها  $T \setminus O$  یک  $1/3$  هسته برای  $S \setminus O$  است؛ پس داریم:

$$KC(S \setminus O, Q -) \leq 3KC(T \setminus O, Q -)$$

حالا توجه کنید  $a$  در آخرین مرحله‌ی الگوریتم  $GMM(T \setminus O, k)$  انتخاب شده است. پس برای هر نقطه  $Q \setminus (T \setminus O) \setminus \{a\}$  چون این نقطه توسط  $GMM$  انتخاب نشده است، پس  $b \in Q -$  نزدیک‌تر از  $a$  است؛ یعنی باید داشته باشیم:  $KC(\{b\}, Q -) \leq KC(\{a\}, Q -)$  این یعنی

$$KC(T \setminus O, Q -) \leq KC(\{a\}, Q -) \leq div(Q)$$

از سه نامساوی قبل حکم نتیجه می‌شود.  $\square$

لم ۸ با احتمال ثابت زمان کوثری  $O((k^2 + \frac{\log n}{r})d \cdot \log k \cdot n^{1/c})$  است.

اثبات. زمان کوثری الگوریتم به ازای هر کوثری سه قسمت دارد. قسمت اول محاسبه  $O(L)$  تابع هش است که  $O(K)$  زمان می‌برد.

$$O(KL) = O((\log_{1/p_r} n) \cdot L) = O(\frac{d}{r} \log n \cdot \frac{\log k}{1 - r/d} \cdot n^{1/c}) = O(\log k \cdot n^{1/c} \cdot \log n \cdot \frac{d}{r})$$

که برای اثبات آن از تقریب  $\log p_r \approx 1 - p_r = \frac{cr}{d}$ ،  $c \geq 1$ ،  $r/d \leq 1/2$  استفاده شده است. دومین قسمت زمان، مربوط به آخرین گام الگوریتم ۵ است که با احتمال ثابت تعداد کل نقاط دورافتاده حداکثر  $3L$  است. پس باید الگوریتم  $GMM$  را برای حداکثر  $O(kL)$  نقطه اجرا کنیم، یعنی  $|T| \leq 3L$ . پس  $GMM$  زمان زیر را می‌برد:

$$O(k^2 L d) = O(d \cdot k^2 \log k \cdot n^{1/c})$$

در نهایت (قسمت سوم زمان الگوریتم)، زمان حلقه داخلی (خطوط ۴ تا ۱۱) الگوریتم ۵، همان طور که قبلاً هم گفته شد می‌تواند طوری پیاده‌سازی شود که کل زمانی که می‌گیرد  $O(kl_i d)$  باشد. پس کل زمان اجرای حلقه  $O(kd \sum l_i) = O(kLd)$  است.  $\square$

## مراجع

- [1] G. Karakostas. A better approximation ratio for the vertex cover problem. *ACM Transactions on Algorithms*, 5(4), 2009.
- [2] I. Dinur, S. Safra. The importance of being biased. In *Proc. 34th ACM Sympos. Theory Comput.*, pages 33–42, 2002.
- [3] S. Khot, O. Regev. Vertex cover might be hard to approximate to within  $2 - \epsilon$ . *Journal of Computer and System Sciences*, 74 (3): 335–349, 2007.