



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده ریاضی و علوم کامپیوتر

پایان نامه کارشناسی ارشد
گرایش سیستم‌های کامپیوتری

کاهش بعد داده‌های بزرگ مقیاس با استفاده از نگاشت
تصادفی

نگارش
سیامک دهد

استاد راهنما
دکتر عادل محمدپور

استاد مشاور
دکتر هادی زارع

دی ۱۳۹۷

صفحه فرم ارزیابی و تصویب پایان نامه - فرم تأیید اعضاء کمیته دفاع

در این صفحه فرم دفاع یا تأیید و تصویب پایان نامه موسوم به فرم کمیته دفاع - موجود در پرونده آموزشی - را قرار دهید.

نکات مهم:

- نگارش پایان نامه/رساله باید به **زبان فارسی** و بر اساس آخرین نسخه دستورالعمل و راهنمای تدوین پایان نامه های دانشگاه صنعتی امیرکبیر باشد.(دستورالعمل و راهنمای حاضر)
- رنگ جلد پایان نامه/رساله چاپی کارشناسی، کارشناسی ارشد و دکترا باید به ترتیب مشکی، طوسی و سفید رنگ باشد.
- چاپ و صحافی پایان نامه/رساله بصورت **پشت و رو(دورو)** بلامانع است و انجام آن توصیه می شود.



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

تعهدنامه اصالت اثر

تاریخ: دی ۱۳۹۷

اینجانب **سیامک دهبید** متعهد می‌شوم که مطالب مندرج در این پایان‌نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آنها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان‌نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است. در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان‌نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان‌نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است. نقل مطالب با ذکر مآخذ بلامانع است.

سیامک دهبید

امضا

تقديم

سپاس‌گزاری

با تشکر از استاد گرامی دکتر محمدپور بابت همراهی و صبر ایشان

سیامک دهبند
دی ۱۳۹۷

چکیده

با ظهور داده‌های بزرگ مقیاس و ناتوانی در نگهداری و پردازش این داده در حافظه، مسئله کاهش بعد اهمیت زیادی پیدا کرده است. یکی از روش‌های کاهش بعد، تصویر تصادفی است که می‌تواند بر روی کلان‌داده‌هایی که بزرگ مقیاس هستند و همچنین بر روی جریان‌های داده، اعمال شود. مبنای این روش ضرب ماتریسی داده‌های اولیه در یک ماتریس تصویرگر است که بعد داده‌های اولیه را کاهش داده ولی اطلاعات آماری مورد نیاز در داده‌های اولیه را با دقت مورد نیاز نگه می‌دارد. روش تصویر تصادفی برای کاهش بعد داده‌های بزرگ مقیاس مزایای متعددی نسبت به روش‌های دیگر کاهش بعد دارد. در این پایان‌نامه این روش برای داده‌های بزرگ مقیاس با دیگر روش‌های کاهش بعد مقایسه شده است. همچنین توانایی این روش برای داده‌هایی با توزیع پایدار غیر نرمال با دیگر روش‌های کاهش بعد مقایسه شده است.

واژه‌های کلیدی:

کاهش بعد، تصویر تصادفی، توزیع پایدار، داده‌های بزرگ مقیاس

فهرست مطالب

عنوان

صفحه

۱	مقدمه	۱
۴	کاهش بعد و داده‌های بزرگ مقیاس	۴
۵	۱-۲ داده‌های حجیم	۵
۵	۱-۱-۲ داده‌های حجیم وب	۵
۷	۲-۱-۲ جریان‌های داده‌ی حجیم	۷
۸	۲-۲ چالش‌های نمونه‌گیری از داده‌های حجیم	۸
۹	۱-۲-۲ مزایای نمونه‌گیری تصادفی مختصات	۹
۹	۲-۲-۲ معایب نمونه‌گیری تصادفی مختصات	۹
۱۰	۳-۲ تصویر تصادفی پایدار	۱۰
۱۱	۴-۲ کاربردها	۱۱
۱۲	۱-۴-۲ کاوش قوانین وابستگی	۱۲
۱۲	۲-۴-۲ وابستگی جفتی همه (فاصله‌ها)	۱۲
۱۳	۳-۴-۲ تخمین فاصله‌ها به طور آنلاین	۱۳
۱۳	۴-۴-۲ بهینه‌سازی درخواست از پایگاه داده	۱۳
۱۴	۵-۴-۲ جستجوی نزدیکترین همسایه از مرتبه‌ی زیر خطی	۱۴
۱۶	۳ تصویر تصادفی پایدار	۱۶
۱۸	۱-۳ مسئله‌ی اصلی در تصویر تصادفی پایدار	۱۸
۱۸	۱-۱-۳ توزیع‌های پایدار	۱۸
۱۹	۲-۱-۳ مسئله برآورد آماری	۱۹
۲۰	۲-۳ تصویر تصادفی نرمال	۲۰
۲۰	۱-۲-۳ مشخصه‌های اصلی	۲۰
۲۳	۳-۳ تصویر تصادفی زیر گوسی و بسیار پراکنده	۲۳
۲۴	۱-۳-۳ تصویر تصادفی زیر گوسی	۲۴
۲۶	۴-۳ تصویر تصادفی کوچی	۲۶

۲۶	۵-۳ تصویر تصادفی α -پایدار
۲۷	۴ کاهش بعد و نحوه‌ی بررسی عملکرد آن
۲۸	۱-۴ PCA و مقایسه با آن
۲۸	۲-۴ برآوردگرای معیارهای وابستگی
۲۸	۳-۴ داده‌های مورد استفاده
۲۸	۴-۴ توضیحات کد
۲۸	۵-۴ شاخص محاسبه عملکرد کاهش بعد و Adjusted Rand Index
۲۹	۵ نتایج
۳۰	منابع و مراجع
۳۵	پیوست
۳۶	واژه‌نامه‌ی فارسی به انگلیسی
۳۸	واژه‌نامه‌ی انگلیسی به فارسی

صفحه	شکل	فهرست اشکال
------	-----	-------------

۱-۲ تصویر تصادفی پایدار $B = A \times R$ ، ماتریس اولیه داده‌ها است. ۱۰

۱-۳ روش تصویر تصافی پایدار ماتریس داده‌ی $A \in \mathbb{R}^{n \times D}$ را در یک ماتریس تصادفی

$R \in \mathbb{R}^{D \times k}$ ضرب می‌کند تا ماتریس تصویر شده‌ی $B = AR \in \mathbb{R}^{n \times k}$ حاصل شود. ۱۷

فهرست جداول

صفحه

جدول

- ۱-۲ تعداد بازدید صفحات برای کلمات با بازخورد بالا و کلمات با بازخورد نادر ۶
- ۲-۲ با افزایش تعداد عبارات در درخواست، باید فرکانس‌های جفت شده کاهش پیدا کنند.
ولی تخمین‌های بیان شده توسط موتورهای جستجو گاهی این موضوع تثبیت شده
را نقض می‌کنند. ۷
- ۳-۲ بازدید صفحات گزارش شده توسط Google برای چهار کلمه و وابستگی‌های دو، سه
و چهارتایی آن‌ها ۱۴

فهرست نمادها

نماد	مفهوم
\mathbb{R}^n	فضای اقلیدسی با بعد n
\mathbb{S}^n	کره n یکه بعدی
M^m	خمینه m -بعدی M
$\mathfrak{X}(M)$	جبر میدان‌های برداری هموار روی M
$\mathfrak{X}^1(M)$	مجموعه میدان‌های برداری هموار یکه روی (M, g)
$\Omega^p(M)$	مجموعه p -فرمی‌های روی خمینه M
Q	اپراتور ریچی
\mathcal{R}	تانسور انحنای ریمان
ric	تانسور ریچی
L	مشتق لی
Φ	۲-فرم اساسی خمینه تماسی
∇	التصاق لوی-چویتای
Δ	لاپلاسین ناهموار
∇^*	عملگر خودالحاق صوری القا شده از التصاق لوی-چویتای
g_s	متر ساساکی
∇	التصاق لوی-چویتای وابسته به متر ساساکی
Δ	عملگر لاپلاس-بلترامی روی p -فرم‌ها

فصل اول

مقدمه

عمومیت پیدا کردن داده‌های حجیم مانند داده‌های حجیم تحت وب و جریان‌های داده بزرگ در کاربردهای جدید، موجب به وجود آمدن فرصت‌های و چالش‌هایی برای مهندسين و دانشمندان شده است. [۴۰] برای مثال، زمانی که ماتریس داده $A \in \mathbb{R}^{n \times D}$ ابعادی در حد وب داشته باشد، عملیات ساده‌ای مانند محاسبه AA^T سخت می‌شود. برای ارائه و نگهداری داده‌های حجیم در حافظه‌ای کوچک و برای استخراج اطلاعات آماری اصلی از مجموعه‌ای از بیانی محدود، روش‌های گوناگونی نمونه‌برداری توسعه یافته است. به طور کلی روش تصویر تصادفی پایدار^۱ برای داده‌های با دم سنگین خیلی خوب کار می‌کند.

روش تصویر تصادفی پایدار، ماتریس داده‌های اولیه $A \in \mathbb{R}^{n \times D}$ را در ماتریس تصادفی $R \in \mathbb{R}^{D \times k}$ ضرب می‌کند و نتیجه ماتریس $B = AR \in \mathbb{R}^{n \times k}$ است. معمولاً درایه‌های ماتریس تصادفی R به صورت i.i.d^۲ از یک توزیع α -پایدار متقارن انتخاب می‌شوند ($0 < \alpha \leq 2$). ما می‌توانیم مشخصه‌های l_α را در A بر اساس B تخمین بزنیم. در مورد حالت l_2 مزیت توزیع تصادفی پایدار توسط لم JL^۳ برجسته شده است. لم بیان می‌دارد که کافی است $k = O(\frac{\log n}{\epsilon^2})$ باشد تا هم فاصله دو به دویی با نرم l_α در A را بتوان با ضریب $1 \pm \epsilon$ از روی ماتریس B تخمین زد. در تز [۴۰] Ping Li لمی مشابه لم JL برای $0 < \alpha < 2$ اثبات شده است. روش تصویر تصادفی پایدار به یک مسئله تخمین آماری کاهش می‌یابد برای تخمین پارامتر مقیاس برای یک توضیح پایدار α متقارن. این مسئله از این جهت مورد توجه قرار می‌گیرد زیرا ما به دنبال برآوردی می‌گردیم که هم از نظر آماری درست باشند و هم از نظر محاسباتی مقرون به صرفه. برآوردگرهای مختلفی را مطالعه و مقایسه کردیم. شامل میانگین حسابی، میانگین هندسی، میانگین هارمونیک، تقسیم توانی^۴ و برآوردگر حداکثر بزرگنمایی.

در این پایان‌نامه ما به بررسی موارد خاصی از تصویر تصادفی پایدار می‌پردازیم. برای نرم l_2 ارتقایی را با استفاده از اطلاعات حاشیه‌ای پیشنهاد می‌کنیم. همچنین برای حالت l_2 می‌توان ماتریس تصویرگر را از یک توزیع زیرگوسی^۵ بسیار کوچکتر به جای توزیع نرمال انتخاب کرد. با در نظر گرفتن محدودیت‌های معقولی می‌توان، از یک توزیع خاص زیرگوسی استفاده کرد. این توزیع شامل $[-1, 0, 1]$ با احتمالات $\{\frac{1}{s}, 1 - \frac{1}{2s}, \frac{1}{s}\}$ با مقادیر بسیار بزرگی برای s (به عبارتی، تصویر تصادف خیلی گسسته^۶) می‌تواند به

¹Stable Random Projection

²Independent and identically distributed random variables

³Johnson-Lindenstrauss

⁴fractional power

⁵sub-Gaussian

⁶Very sparse random projections

خوبی تصویر تصادفی نرمال عمل کند. برای حالت نرم l_1 به عبارتی دیگر تصویر تصادفی کوچی^۷ انجام تخمین کاری نسبتاً جذاب است. برای مثال، محاسبه برآوردگر بیشینه درستنمایی MLE در این حالت از لحاظ محاسباتی ممکن است. و یک توزیع معکوس گاوسی^۸ برای مدل سازی دقیق توزیع MLE بیان شده است.

روش تصویر تصادفی از پراکندگی داده ها استفاده ای نمی کند. در حالی که داده های بزرگ مقیاس معمولاً بسیار پراکنده هستند. از روش تصویر تصادفی می توان برای حل مسائل بزرگ مقیاس در علوم و مهندسی در موتورهای جستجو و سیستم های اخذ داده، پایگاه های داده، سیستم های جریان داده جدید، جبر خطی عددی و بسیاری از کارهای یادگیری ماشین و داده کاوی که شامل محاسبه حجم فاصله ها است، استفاده کرد.

⁷Cauchy random projections

⁸inverse Gaussian

فصل دوم

کاهش بعد و داده‌های بزرگ مقیاس

۱-۲ داده‌های حجیم

عبارات زیر از سایت *Information Week* نقل قول شده‌اند^۱:

- مقدار داده‌ای که توسط کسب و کارها ذخیره می‌شود تقریباً هر ۱۲ تا ۱۸ ماه دو برابر می‌شود.
 - پایگاه داده‌ها بیشتر هم زمان شده‌اند. فروشگاه‌های زنجیره‌ای Wall-Marat داده‌های فروش را هر ساعت به روز می‌کند.
 - اضافه شدن یک میلیون خط داده اجازه جستجوهای پیچیده‌تری را می‌دهد. شرکت EBay به کارمندان اجازه می‌دهد برای بدست آوردن درکی عمیق‌تر در خصوص رفتار مشتریان در میان داده‌های حراج در بازه‌های زمانی کوتاه جستجو کنند.
 - بزرگترین پایگاه داده‌ها توسط، مرکز شتاب‌دهنده خطی استاندارد، مرکز تحقیقات ناسا، آژانس امنیت ملی و ... در ابعادی در محدوده‌ی پتابایت (هزار ترابایت 10^{15} بایت)، اداره می‌شوند.
- پدیده نو ظهور مجموعه داده‌ای حجیم، چالش‌های محاسباتی در بسیاری کاربردهای علمی و تجاری به وجود آورده است. شامل اختریف‌یک، بیوتکنولوژی، جمعیت‌شناسی^۲، مالی، سیستم‌های اطلاعات جغرافیایی، دولت، دارو، ارتباطات از راه دور، محیط زیست و اینترنت.

۱-۱-۲ داده‌های حجیم وب

وب چقدر بزرگ است؟ **جدول ۱-۲** نشان‌دهنده تعداد بازدید صفحات در موتورهای جستجوی امروزی است. به طور تخمینی حدود $D = 10^{10}$ صفحه‌ی وب را می‌توان بر اساس بازدید دو واژه‌ی بسیار پر کاربرد «A» و «THE» تخمین زد. **جدول ۱-۲** همچنین نشان می‌دهد که حتی کلماتی که به ندرت کاربرد دارند هم تعداد زیادی بازدید دارند.

کلماتی با بازخورد معمولی چه میزان بازدید دارند؟ برای جواب این سوال ما به طور تصادفی ۱۵ صفحه از لغتنامه‌ی آموزشی انتخاب می‌کنیم. [۳۱] (لغتنامه‌ای با ۵۷،۱۰۰ کلمه) و اولین کلمه در هر صفحه را مد نظر قرار می‌دهیم. میانه‌ی آماری بر اساس جستجوگر گوگل ۱۰ میلیون صفحه برای کلمه است.

زبان انگلیسی چند کلمه دارد؟ در اینجا عبارتی را از AskOxford.com نقل قول می‌کنیم:

^۱<http://www.informationweek.com/news/showArticle.jhtml?articleID=175801775>

^۲demographics

Query	Google	Bing
A	25,270,000,000	175,000,000
The	25,270,000,000	101,000,000
Kalevala	7,440,000	939,000
Griseofulvin	1,163,000	332,000
Saccade	1,030,000	388,000

جدول ۱-۲: تعداد بازدید صفحات برای کلمات با بازخورد بالا و کلمات با بازخورد نادر

« این بیان میدارد که حداقل یک چهارم میلیون واژه‌ی انگلیسی مستقل وجود دارد. به جز افعال صرفی و کلمات فنی و ناحیه‌ای که توسط OED^۳ تحت پوشش قرار نمی‌گیرند یا کلماتی که هنوز به لغتنامه‌های منتشر شده اضافه نشده‌اند. در صورتی که این موارد هم در نظر گرفته شوند تعداد لغات در حدود سه چهارم میلیون لغت خواهد بود »

بنابراین اگر یک ماتریس «عبارت به سند» $A \in \mathbb{R}^{n \times D}$ در نظر بگیریم. در ابعاد وب این ماتریس در ابعاد $n \approx 10^6$ و $D \approx 10^{10}$ بزرگ خواهد شد. در اینجا عدد (i, j) در A تعداد ظهور واژه i در سند j را نشان می‌دهد.

کارکردن با ماتریسی در این ابعاد بزرگ چالش برانگیز است. برای مثال، شاخص LSI^۴ [۲۳] و یک مدل موضوعی فراگیر، از SVD^۵ بر روی ماتریس عبارت به سند استفاده می‌کند. که انجام این عملیات در ابعاد وب قطعاً غیرممکن است.

یک مشکل اصلی در قبال مجموعه داده‌های سنگین، حافظه کامپیوتر است. به این دلیل که ابعاد و سرعت حافظه فیزیکی بسیار رشد کمتری در مقایسه با پردازنده‌ها (CPU) دارد. این پدیده به عنوان دیوار حافظه شناخته می‌شود [۴۳، ۴۸]. برای مثال، هر چند ممکن است تمامی رخدادهای همزمان دوتایی از پیش محاسبه شوند، ولی نگهداری این حجم از داده در حافظه غیر ممکن است. علاوه بر این، گاهی اوقات تخصیص‌هایی با بیش از دو عامل هم اهمیت پیدا می‌کنند زیرا درخواست‌ها ممکن است شامل بیش از دو واژه هم باشند. یک راه حل ممکن این است که یک «نمونه» از A نگهداری شود و همزمانی‌ها بر اساس این نمونه در حین کار تخمین زده شوند. ما حدس می‌زنیم که این روش توسط موتورهای جستجوی امروزی مورد استفاده قرار می‌گیرد، هر چند که روش واقعی قطعاً جزو اسرار تجاری آن‌ها است.

هر چند که انتظار می‌رود تخمین‌ها سازگار باشند و فرکانس‌های جفت شده باید با افزایش عبارت به

³Oxford English Dictionary

⁴latent semantic indexing

⁵singular value decomposition

درخواست، کاهش پیدا کنند. **جدول ۲-۲** نشان می‌دهد که تخمین‌های بیان شده با موتورهای جستجوی فعلی، همیشه سازگار نیستند.

Query	Hits(Bing)	Hits(Google)
America	150,731,182	393,000,000
America & China	15,240,116	66,000,000
America & China & Britain	235,111	6,090,000
America & CHina & Britain & Japan	154,444	23,300,000

جدول ۲-۲: با افزایش تعداد عبارات در درخواست، باید فرکانس‌های جفت شده کاهش پیدا کنند. ولی تخمین‌های بیان شده توسط موتورهای جستجو گاهی این موضوع تثبیت شده را نقض می‌کنند.

با اینکه، تعداد کل واژه‌های انگلیسی (که به‌طور صحیح نوشته شده‌اند) هم اکنون شگفت‌آور است، در بسیاری کاربردهای متن کاوی، ما باید با ابعاد بسیار بزرگتری سر و کار داشته باشیم. در حالی که یک سند ممکن است بیانگر برداری از تک واژه‌ها باشد (به عبارت دیگر، مدل کیسه لغات ^۶). معمولاً بهتر است سند به عنوان یک بردار از لغات به صورت ۱ پیوسته ^۷ [۱۵] بیان شود. برای مثال، با استفاده از مدل ۳ پیوسته، جمله‌ی "It is a nice day" به مجموعه‌ی زیر تجزیه می‌شود. "a", "is a nice", "it is a" {"it is a", "is a nice", "a"}
 nice day" این مدل به طور چشمگیری ابعاد داده‌ها را افزایش می‌دهد. به خاطر اینکه، اگر مجموعه‌ی 10^6 تک لغت انگلیسی موجود داشته باشد. مدل ۳ پیوسته تعداد ابعاد را از 10^6 به 10^{18} افزایش می‌دهد.

۲-۱-۲ جریان‌های داده‌ی حجیم

در بسیاری کاربردهای جدید پردازش داده، جریان‌های داده‌ی حجیم نقش بنیادی دارند. جریان‌های داده‌ای که از روترهای اینترنت، سوئیچ‌های تلفن، رصد امسفر، شبکه‌های سنسور، شرایط ترافیکی بزرگرایی، داده‌های مالی و غیره [۳، ۴۴، ۲۲، ۹، ۳۲، ۳۷، ۳۰] حاصل می‌شوند.

برخلاف پایگاه داده‌های سنتی، معمول نیست که جریان‌های داده‌ی حجیم (که با سرعت زیادی منتقل می‌شوند) در جای نگهداری شوند. بنابراین پردازش معمولاً به طور همزمان انجام می‌شوند. برای مثال، گاهی اوقات «رصد تصویری» داده‌ها با رصد تغییرات زمانی برخی آماره‌ها کفایت می‌کند. برای مثال آماره‌های نظیر: مجموع، تعداد آیتم‌های مجزا، برخی نرم‌های l_α . در برخی کاربردها (برای مثال، طبقه‌بندی صدا/محتوا و جدا سازی) نیاز است یک مدل یادگیری آماری برای کلاسه‌بندی ^۸ یا خوشه‌بندی

^۶bag-of-words

^۷l-shingles

^۸Classification

^۹ جریان داده‌های حجیم توین شود. ولی معمولاً فقط می‌توانیم یک‌بار داده‌ها را مورد بررسی قرار دهیم. یک خاصیت مهم جریان‌های داده‌ای این است که دینامیک هستند. به عنوان یک مدل محبوب، جریان u شامل ورودی‌های (i, u_i) است که $i = 1 to D$. برای مثال، $D = 2^{64}$ زمانی که جریان بیان‌گر IP آدرس‌ها است. ^{۱۰} ورودی‌ها ممکن است به هر ترتیبی باشند و ممکن است مرتباً به روز شوند. ذات دینامیک جریان داده‌های حجیم فرآیند نمونه‌گیری را بسیار چالش‌برانگیزتر از زمانی می‌کند که با داده‌های ایستا سر و کار داریم.

۲-۲ چالش‌های نمونه‌گیری از داده‌های حجیم

در حالی که مسائل جذاب و چالش‌برانگیزی با ورود داده‌های حجیم شکل گرفته‌اند، این پایان‌نامه بر روی توسعه‌ی روش‌های نمونه‌گیری برای محاسبه فاصله در داده‌هایی با ابعاد بسیار بالا با استفاده از حافظه محدود تمرکز دارد.

در کاربردهای مدل‌سازی آماری و یادگیری ماشین، در اغلب موارد به جای داده‌های اصلی به فاصله، به خصوص فاصله‌ی جفتی نیاز داریم. برای مثال، محاسبه ماتریس گرام AA^T ^{۱۱} در آمار و یادگیری ماشین معمول است. AA^T بیانگر همه‌ی ضرب‌های داخلی دوتایی در ماتریس داده‌ی A است. دو داده‌ی $u_1, u_2 \in \mathbb{R}^D$ داده شده‌اند. ضرب داخلی آن‌ها (که با a نمایش داده می‌شود) و l_α (که با $d_{(\alpha)}$ نمایش داده می‌شوند با عبارات زیر تعریف می‌شوند: ^{۱۲}

$$a = u_1^T u_2 = \sum_{i=1}^D u_{1,i} u_{2,i} \quad (1-2)$$

$$d_{(\alpha)} = \sum_{i=1}^D |u_1 - u_2|^\alpha \quad (2-2)$$

^۹Clustering

^{۱۰} هرچند ما بیشتر اوقات تعداد دقیق ابعاد (D) یک جریان داده را نمی‌دانیم ولی در بیشتر کاربردها کافی است حد بالایی محافظه‌کارانه‌ای را در نظر بگیریم. برای مثال $D = 2^{64}$ زمانی که جریان بیانگر IP های ورودی است. همچنین این یکی از دلایلی است که داده‌ها بسیار پراکنده هستند. به این نکته توجه داشته باشید که ابعاد بسیار بزرگ تأثیری در محاسبه‌ی فاصله‌ها و نمونه‌گیری طی الگوریتم‌های معرفی شده در این پایان‌نامه ندارد.

^{۱۱}Gram matrix

^{۱۲} ما فاصله l_α را به صورت $d_{(\alpha)} = \sum_{i=1}^D |u_1 - u_2|^\alpha$ تعریف کرده‌ایم. به جای اینکه به شکل $(\sum_{i=1}^D |u_1 - u_2|^\alpha)^{1/\alpha}$ تعریف کنیم. زیرا شکل اول در کاربردهای عملی عمومیت بیشتری دارد. برای مثال، l_2 ، در ادبیات معمولاً به شکل توان دو l_2 بیان می‌شود. $\sum_{i=1}^D |u_1 - u_2|^2$ به جای $(\sum_{i=1}^D |u_1 - u_2|^2)^{1/2}$. در این پایان‌نامه، ما برای سادگی $\sum_{i=1}^D |u_1 - u_2|^2$ را «فاصله l_2 » بیان می‌کنیم به جای «مربع فاصله‌ی l_2 ».

به این نکته توجه داشته باشید که هم ضرب داخلی و هم فاصله به شکل جمع D جمله تعریف می‌شوند. بنابراین، زمانی که داده‌ها به اندازه‌ای حجیم هستند که نمی‌توان به طور کارا آن‌ها را مدیریت کرد، نمونه‌گیری خیلی عادی به نظر می‌رسد تا بتوان با انتخاب تصادفی k عضو از D جمله تخمینی از مجموع به دست آوریم (با ضریب مقیاس $\frac{D}{k}$). در خصوص ماتریس داده‌ی $A \in \mathbb{R}^{n \times D}$ انتخاب تصادفی مختصات ^{۱۳}، k ستون را از ماتریس داده به طور یکنواخت و تصادفی انتخاب می‌کند.

نمونه‌گیری از این جهت سودمند است که هم سایکل‌های کاری CPU را کاهش می‌دهد و هم در حافظه صرفه‌جویی می‌کند. در کاربردهای جدید، در اغلب موارد صرفه‌جویی در حافظه از اهمیت بیشتری برخوردار است. در نیم قرن گذشته گلوگاه محاسباتی حافظه بوده است، نه پردازشگر. سرعت پردازشگرها با نرخ تقریبی ۷۵ درصد در سال رو به افزایش است. در حالی که سرعت حافظه تقریباً سالی ۷ درصد افزایش می‌یابد [۴۳]. این پدیده به عنوان «دیوار حافظه» ^{۱۴} شناخته می‌شود [۴۳، ۴۸]. بنابراین در کاربردهایی که شامل مجموعه داده‌های حجیم می‌شوند، بحرانی‌ترین کار بیان کردن داده‌ها است. برای مثال، از طریق نمونه‌گیری با فرمی فشرده برای قرارگیری در ابعاد حافظه در دسترس.

۱-۲-۲ مزایای نمونه‌گیری تصادفی مختصات

نمونه‌گیری تصادفی مختصات به دو دلیلی معمولاً انتخاب پیش‌فرض است.

- **سادگی** این روش از لحاظ زمانی تنها از مرتبه $O(nk)$ برای نمونه‌گیری k ستون از $A \in \mathbb{R}^{n \times D}$ طول می‌کشد.
- **انعطاف‌پذیری** یک مجموعه نمونه را می‌توان برای تخمین بسیاری از شاخص‌های آماری استفاده کرد. شامل: ضرب داخلی، فاصله l_α (برای هر مقداری از α)

۲-۲-۲ معایب نمونه‌گیری تصادفی مختصات

با این حال نمونه‌گیری تصادفی مختصات دو ایراد اساسی دارد.

- معمولاً دقیق نیست زیرا مقادیری با مقدار زیاد محتمل است که گم شوند. مخصوصاً زمانی که داده‌ها دم سنگینی داشته باشند. داده‌های بزرگ مقیاس دنیای واقعی (مخصوصاً داده‌های مربوط به اینترنت) همیشه دم سنگین هستند و از قاعده توانی پیروی می‌کنند. [۳۹، ۲۰، ۲۵، ۴۵]

¹³Random coordinate sampling

¹⁴Memory wall

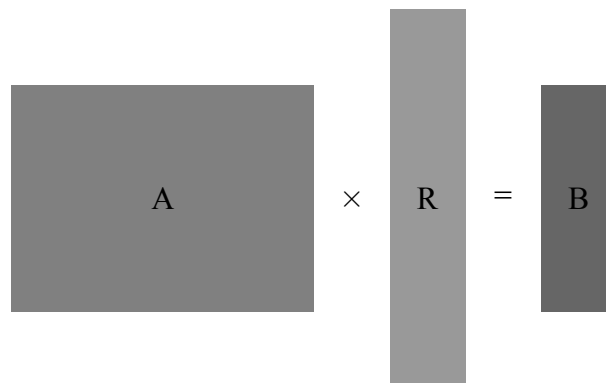
که فاصله l_2 یا ضرب داخلی را تخمین می‌زنیم. واریانس تخمین‌ها بر اساس ممان چهارم داده‌ها تعیین می‌شود. در حالی که در داده‌های دم سنگین، گاهی اوقات حتی ممان اول هم معنی‌دار نیست (محدود نیست) [۴۵].

• این روش داده‌های پراکنده را به خوبی مدیریت نمی‌کند. بسیاری از داده‌های بزرگ مقیاس به شدت پراکنده هستند، به عنوان مثال، داده‌های متنی [۲۴] و داده‌های بر اساس بازار [۴۶، ۴]. به جز برخی واژه‌های کاربردی مانند "A" و "The" بیشتر لغات با نسبت بسیار کمی در مستندات ظاهر می‌شوند ($< 1\%$) اگر ما داده‌ها را با در نظر گرفتن تعدادی از ستون‌های ثابت نمونه‌گیری کنیم. خیلی محتمل است که بیشتر داده‌های (مقادیر غیر صفر) را از دست بدهیم. به خصوص موارد جذابی که دو مقدار با هم غیر صفر شده‌اند.

در این پایان‌نامه ما روش تصویر تصادفی را مورد بررسی قرار می‌دهیم و نشان خواهیم داد که این روش به خوبی قابلیت مدیریت داده‌های دم‌سنگین را دارد.

۳-۲ تصویر تصادفی پایدار

تصویر شکل ۱-۲، ایده تصویر تصادفی را نشان می‌دهد. ایده اصلی تصویر تصادفی ضرب ماتریس داده‌ی $A \in \mathbb{R}^{n \times D}$ در ماتریس تصادفی $R \in \mathbb{R}^{D \times k}$ ($k \ll D$) است. که حاصل ماتریس تصویر شده‌ی $B = A \times R \in \mathbb{R}^{n \times k}$ است. B بسیار کوچکتر از A است و بنابراین به راحتی قابل ذخیره‌سازی است. (برای مثال: برای حافظه‌های فیزیکی به اندازه‌ی کافی کوچک است)



شکل ۱-۲: تصویر تصادفی پایدار $B = A \times R$ ، A ماتریس اولیه داده‌ها است.

ماتریس تصویرگر $R \in \mathbb{R}^{D \times k}$ معمولاً از داریه‌های مستقل هم توزیع (i.i.d) یک توزیع متقارن α -پایدار پر شده است. [۴۹] (بنابراین نام این روش «تصویر تصادفی پایدار» است). بر اساس مشخصات

توزیع‌های α -پایدار، داده‌های تصویر شده هم از توزیع α -پایدار پیروی می‌کنند. که بر اساس آن‌ها شاخص‌های l_α و فاصله دودویی l_α در A تخمین زده می‌شوند و می‌توانیم داده‌های اصلی را دور بریزیم. موفقیت تصویر پایدار تصادفی توسط لم Johnson-Lindenstrauss (JL) [۳۵] برای کاهش بعد در l_2 نشان داده شده است. لم JL بیان می‌کند: رعایت $k = O\left(\frac{\log n}{\epsilon^2}\right)$ تضمین می‌کند هر فاصله l_2 میان n نقطه در هر تعداد بعدی با دقت $1 \pm \epsilon$ با احتمال بالایی تخمین زده شود. (k در اینجا بیانگر تعداد ابعاد کاهش یافته است)

با این حال لم JL برای نرم‌های فاصله با α کوچکتر از ۲ ($\alpha < 2$) صادق نیست. در صورتی که لازم باشد از برآوردگرهایی استفاده کنیم که متریک باشند (در نامساوی مثلثی صدق کنند). به این نتیجه «عدم امکان»^{۱۵} گفته می‌شود. [۱۶، ۳۸، ۱۴] خوشبختانه شامل برآوردگرهایی که متریک نیستند نمی‌شود. در این پایان‌نامه ما در مورد برآوردگرهای کوناگونی که متریک نیستند صحبت خواهیم کرد. شامل: میانگین هندسی^{۱۶}، میانگین هارمونیک^{۱۷}، نسبت توانی^{۱۸} و همچنین حداکثر بزرگنمایی.

۴-۲ کاربردها

علاقه‌ی زیادی به تکنیک‌های نمونه برداری وجود دارد که در کاربردهای زیادی مورد استفاده قرار می‌گیرند. مانند: قانون وابستگی^{۱۹} [۱۱، ۱۲]، خوشه‌بندی، بهینه‌سازی درخواست^{۲۰} [۴۲، ۱۸]، تشخیص تکراری^{۲۱} [۱۵، ۱۰] و بسیاری موارد دیگر. روش‌های نمونه بردار هر چه بیشتر و بیشتر برای مجموعه‌های بزرگتر اهمیت پیدا می‌کنند.

طرح برودر^{۲۲} [۱۵] در ابتدا برای تشخیص صفحات وب تکراری معرفی شد. URL‌های زیادی به HTML‌های مشابه (یا تقریباً مشابه) اشاره می‌کنند. جواب‌های تخمین زده شده به اندازه‌ی کافی خوب بودند. نیازی نبود تا همه تکراری‌ها پیدا شوند ولی کاربردی بود که تعداد زیادی از آن‌ها پیدا شوند، بدون اینکه بیش از ارزش آن از توان محاسباتی استفاده شود.

در کاربردهای بازایی اطلاعات (IR)^{۲۳} معمولاً گلوگاه حافظه‌ی فیزیکی است. زیرا مجموعه‌ی

¹⁵Impossibility

¹⁶Geometric mean

¹⁷Harmonic mean

¹⁸Fractional power

¹⁹Association rules

²⁰Query optimization

²¹Duplicate detection

²²Broder's sketch

²³information retrieval

وب برای حافظه (RAM) بسیار بزرگ است و از طرفی ما می‌خواهیم زمان گشتن به دنبال داده^{۲۴} بر روی دیسک را کمینه کنیم. زیرا زمان پاسخ به یک درخواست کلیدی است [۱۳]. به عنوان یک وسیله صرفه‌جویی در فضا، کاهش بعد یک ارائه فشرده از داده‌ها فراهم می‌کند که برای تولید جواب‌های تخمینی در حافظه فیزیکی مورد استفاده قرار می‌گیرند.

ما به بازدید صفحات وب اشاره کردیم. اگر ما یک عبارت جستجوی دو کلمه‌ای داشته باشیم، می‌خواهیم بدانیم چه تعداد از صفحات هر دو کلمه را دارند. فرض می‌کنیم محاسبه‌ی از قبل و نگهداری بازدید صفحات غیر ممکن باشد. حداقل نه برای کلماتی که تکرار زیادی ندارند و سری‌های چند کلمه‌ای. مرسوم است که در بازیابی اطلاعات با یک ماتریس بزرگ عبارت به ازای سند شروع کنیم که در آن مقادیر ورودی نشان‌دهنده‌ی وجود عبارت در متن است. بنا به کاربردهای خاص می‌توانیم بک اندیس معکوس^{۲۵} بسازیم و کلیتی از عبارات (برای تخمین ارتباط لغات) یا اسناد (برای تخمین شباهت اسناد) نگهداری کنیم.

۱-۴-۲ کاوش قوانین وابستگی

تحلیل‌های مبتنی بر بازار و قوانین وابستگی [۵، ۶، ۷] ابزارهای مناسبی برای کاوش پایگاه داده‌های تجاری هستند. پایگاه داده‌های تجاری دارند روز به روز بزرگتر و گسسته‌تر می‌شوند. [۴، ۴۶] الگوریتم‌های مختلف نمونه‌برداری پیشنهاد شده است. نمونه برداری این امکان را فراهم می‌کند تا قواعد تخصیص را به صورت آنلاین برآورد کنیم. که می‌تواند مزایایی در کاربردهای خاص داشته باشد.

۲-۴-۲ وابستگی جفتی همه (فاصله‌ها)

در کاربردهای مختلفی شامل کلاسه‌بندی بر مبنای فاصله یا خوشه‌بندی و مدل‌سازی زبان با bi-gram^{۲۵} ما نیازمند محاسبه‌ی همه‌ی جفت تخصیص‌ها (یا فاصله‌ها) هستیم. ماتریس داده‌ی A شامل n سطر و D ستون داده شده است. محاسبه‌ی مستقیم AA^T ، $O(n^2D)$ هزینه بر است. یا به طور بهینه‌تر $O(n^2\bar{f})$ که \bar{f} تعداد میانگین مقادیر غیر صفر میان تمام سطرها A است. محاسبه مستقیم می‌تواند به شدت زمان‌بر باشد. همچنین، به طور خاص زمانی که ماتریس داده آنقدر بزرگ است که در حافظه فیزیکی جا نمی‌شود. محاسبه به طور خاص بسیار ناکارآمد خواهد بود.

^{۲۴}inverted index

^{۲۵}litez48

۳-۴-۲ تخمین فاصله‌ها به طور آنلاین

در حالی که ماتریس داده‌ی اولیه $A \in \mathbb{R}^{n \times D}$ ممکن است برای حافظه‌ی فیزیکی بسیار بزرگ باشد، نگهداری^{۲۶} همه فاصله‌های جفتی و وابستگی‌ها در A ، $O(n^2)$ فضا مصرف می‌کند. که می‌تواند برای حافظه‌ی فیزیکی بسیار بزرگتر باشد. در این میان وابستگی‌های چندتایی را کنار می‌گذاریم. در بسیاری از کاربردها نظیر یادگیری برخط، سیستم‌های توصیه آنلاین، تحلیل‌های بازار برخط و موتورهای جستجو، بهتر است که نمونه‌ها (sketches) در حافظه نگهداری شوند و همه‌ی فاصله‌ها به طور آنلاین، زمانی که مورد نیاز باشد، محاسبه شوند.

۴-۴-۲ بهینه‌سازی درخواست از پایگاه داده

در پایگاه داده‌ها یک وظیفه‌ی بسیار مهم تخمین join‌های چندراهی است، که تاثیر زیادی بر روی کارایی سیستم دارد [۲۹]. بر اساس تخمین دوراهی، سه‌راهی و حتی join‌هایی از مرتبه‌ی بالاتر، بهینه‌گرهای درخواست یک نقشه برای کمینه کردن تابع هزینه می‌سازند (برای مثال، نوشتن‌های میانی^{۲۷}). بهینه بودن اهمیت بسیاری دارد زیرا مثلاً نمی‌خواهیم زمان بیشتری برای بهینه‌سازی نقشه نسبت به زمان اجرای آن تلف کنیم.

ما از مثال Governor برای نمایش کاربرد تخمین دو و چند راهه برای بهینه کردن درخواست استفاده می‌کنیم.

جدول ۳-۲ بازدید صفحات را برای چهار کلمه و ترکیبات دو، سه، چهارتایی آن‌ها نشان می‌دهد. فرض

کنیم بهینه‌ساز قصد استخراج نقشه برای درخواست: "Governor, Schwarzenegger, Terminator, Austria" را داشته باشد. راه حل استاندارد این است که با عبارات با کمترین فراوانی شروع کند: ("Schwarzenegger" \cap "Governor") \cap "Austria" این نقشه 579,100 نوشتن میانی بعد از اولین و دومین join‌ها دارد. یک بهینه‌سازی می‌تواند ("Terminator" \cap "Austria") \cap ("Schwarzenegger") \cap "Governor" باشد که 579,100 را به 136,000 کاهش می‌دهد.

²⁶Materializing

²⁷Intermediate writes

	Query	Hits(Google)
One-way	Austria	88,200,000
	Governor	37,300,000
	Schwarzenegger	4,030,000
	Terminator	3,480,000
Two-way	Governor & Schwarzenegger	1,220,000
	Governor & Austria	708,000
	Schwarzenegger & Terminator	504,000
	Terminator & Austria	171,000
	Governor & Terminator	132,000
	Schwarzenegger & Austria	120,000
Tree-way	Governor & Schwarzenegger & Terminator	75,100
	Governor & Schwarzenegger & Austria	46,100
	Schwarzenegger & Terminator & Austria	16,000
	Governor & Terminator & Austria	11,500
Four-way	Governor & Schwarzenegger & Terminator & Austria	6,930

جدول ۲-۳: بازدید صفحات گزارش شده توسط Google برای چهار کلمه و وابستگی‌های دو، سه و چهارتایی آن‌ها

۵-۴-۲ جستجوی نزدیکترین همسایه از مرتبه‌ی زیر خطی

محاسبه‌ی نزدیکترین همسایه در بسیاری کاربردها از اهمیت زیادی برخوردار است. با این حال، به دلیل «نفرین ابعاد»^{۲۸} راه حل فعلی برای پیدا کردن بهینه‌ی نزدیکترین همسایه‌ها (حتی به طور تقریبی) اصلاً رضایت بخش نیست. [۲۶، ۳۴]

به دلیل ملاحظات محاسباتی، دو شکل اصلی در جستجوی نزدیکترین همسایه‌ها وجود دارد. اول اینکه ماتریس اصلی داده‌ها $A \in \mathbb{R}^{n \times D}$ ممکن است برای حافظه فیزیکی بسیار بزرگ باشد ولی اسکن کردن دیسک‌های سخت برای پیدا کردن نزدیکترین همسایه‌ها می‌تواند خیلی کند باشد. دوماً، پیدا کردن نزدیکترین همسایه‌های یک داده ممکن است $O(nD)$ هزینه بر باشد که می‌تواند به شدت زمان بر شود.

با این حال، روس کاهش ابعادی در این پایان‌نامه می‌تواند در حافظه صرفه‌جویی کند و سرعت محاسبات را افزایش دهد. برای مثال: وقتی ماتریس داده‌ی اولیه A به ماتریس داده‌ی $B \in \mathbb{R}^{n \times k}$ کاهش می‌یابد. با این حال، $O(nk)$ و معمولاً این درخواست وجود دارد که هزینه‌ی محاسباتی از $O(n)$ به $O(n^\gamma)$ برای $\gamma < 1$ کاهش پیدا کند، حداقل برای کاربردهای خاص.

دو گروه اصلی الگوریتم‌های زیر خطی برای محاسبه عبارتند از KD-Trees (و انواع آن) [۲۷، ۲۸]

²⁸Curse of dimensionality

و Locality-Sensitive Hashing (LSH) [۸، ۲۱، ۳۴] این الگوریتم‌ها معمولاً با یک فضای متریک کار می‌کنند (که در آن نامساوی مثلثی برقرار است). برای مثال، فضای l_α زمانی که $\alpha \geq 1$ باشد یک متریک است. زمانی که به دنبال نزدیکترین همسایه‌ها در l_α ($\alpha > 1$) می‌گردیم، می‌توانیم (نسبتاً به سادگی) فضای جستجو را به طور کاملاً اساسی با استفاده از نامساوی مثلثی کاهش دهیم. به عبارت دیگر، نیازی نیست که همه n نقطه داده‌ها را مورد بررسی قرار دهیم.

در داده‌هایی با ابعاد بسیار بزرگ، الگوریتم‌های زیر خطی موجود شامل KD-trees و LSH، عملکرد رضایت بخشی ندارند. وقتی حافظه‌ی فیزیکی (به جای CPU) گلوگاه باشد^{۲۹}، یکی از مشکلات اصلی این است که این الگوریتم‌ها برای کاهش هزینه‌ی محاسباتی به حافظه‌ی ابر خطی^{۳۰} نیاز دارند که می‌تواند مشکل ساز باشد. [۳۴] به طرح کلی برای LSH توجه کنید که ترکیبی از هش^{۳۱} و تصویر تصادفی است. متأسفانه این طرح به دلیل هزینه‌ی زیاد پیش پردازش غیر کاربردی است. [۳۴]

در این پایان‌نامه، موفقیت اصلی کاهش بعد داده $A \in \mathbb{R}^{n \times D}$ به $B \in \mathbb{R}^{n \times k}$ و تامین برآوردهای دقیق برای استخراج فاصله‌های اولیه در A بر اساس B است. در حالی که سناریوهای مهمی وجود دارند که در آن‌ها نتایج ما رضایت بخش هستند، توسعه‌ی یک الگوریتم زیر-خطی برای تخمین نزدیکترین همسایه‌ها، بر اساس الگوریتم ما یک ایده جذاب برای تحقیقات آینده است. یک مانع اصلی در این راه این است که بیشتر برآوردهای ما غیر متریک هستند و بنابراین طراحی یک الگوریتم هوشمند و تحلیل‌های تئوری ممکن است سخت باشد، با این حال غیر ممکن نیست.

²⁹Memory wall³⁰Super-linear memory³¹Hash

فصل سوم

تصویر تصادفی پایدار

روش تصویر تصادفی پایدار^۱ [۴۱، ۴۷، ۹، ۳۳، ۳۲، ۳۶] یک روش پرکاربرد در داده‌کاوی و یادگیری ماشین است. با این روش به طور کار $l_\alpha (0 < \alpha \leq 2)$ فاصله در داده‌های حجیم (برای مثال: وب یا جریان‌های داده‌ی حجیم) محاسبه می‌شود. در این روش حافظه‌ی کمی استفاده شده و فقط یک بار پایش داده‌ها کافی است.

$$A \times R = B$$

شکل ۳-۱: روش تصویر تصادفی پایدار ماتریس داده‌ی $A \in \mathbb{R}^{n \times D}$ را در یک ماتریس تصادفی $R \in \mathbb{R}^{D \times k}$ ضرب می‌کند تا ماتریس تصویر شده‌ی $B = AR \in \mathbb{R}^{n \times k}$ حاصل شود.

همانطور که در شکل ۳-۱ می‌بینید، ایده تصویر تصادفی پایدار، ضرب ماتریس داده‌ها $A \in \mathbb{R}^{n \times D}$ در ماتریس تصادفی $R \in \mathbb{R}^{D \times k} (k \ll D)$ است که حاصل یک ماتریس تصویر شده‌ی $B \in \mathbb{R}^{n \times k}$ است. درایه‌های ماتریس تصادفی \mathbf{R} به طور i.i.d. (مستقل و هم توزیع)^۲ از یک توزیع α -پایدار^۳ حاصل می‌شوند. به همین دلیل به این روش «تصویر تصادفی پایدار» گفته می‌شود. به این نکته توجه کنید که توزیع ۲-پایدار معادل توزیع نرمال و توزیع پایدار ۱-پایدار معادل کوچی^۴ است. حالت خاص تصویر تصادفی نرمال (به عبارت دیگر $\alpha = 2$) نسبتاً به خوبی مورد بررسی قرار گرفته است. به رساله [۴۷] مراجعه کنید. بنابراین، بخش اعظم این پایان‌نامه به تصویر تصادفی پایدار $\alpha < 2$ اختصاص یافته است.

پس از مروری بر حالت کلی تصویر تصادفی پایدار $0 < \alpha \leq 2$ ، جزئیات بیشتری در خصوص حالت l_2 مورد بررسی قرار می‌گیرد. سپس ارتقاء روش با استفاده از اطلاعات حاشیه‌ای^۵ بررسی می‌شود. در ادامه، تصویر تصادفی نرمال ساده‌سازی می‌شود. این کار با نمونه‌برداری R از حالت توزیع گسسته‌ی

¹Stable Random Projections

²Independent and Identically distributed

³ α -stable distribution

⁴Cauchy

⁵Marginal information

سه نقطه‌ای $[-1, 0, 1]$ انجام می‌شود. این حالت، یک حالت خاص توزیع‌های زیرگوسی^۶ است. سپس نرم l_1 ^۷ مورد بررسی قرار گرفته و در ادامه حالت کلی $0 < \alpha \leq 2$ مورد بحث قرار می‌گیرد.

۱-۳ مسئله اصلی در تصویر تصادفی پایدار

مسئله اصلی تصویر تصادفی پایدار یک مسئله برآورد آماری است. همانطور که بیان شد، ماتریس داده‌ی $A \in \mathbb{R}^{n \times D}$ را در ماتریس تصادفی $R \in \mathbb{R}^{D \times k}$ ضرب می‌کنیم تا ماتریس بسیار کوچکتر $B = A \times R \in \mathbb{R}^{n \times k}$ را بدست بیاوریم. هدف این است که مشخصات آماری A بر اساس ماتریس B استنتاج شوند. (شامل نرم و فاصله)

بدون از دست دادن کلیت، ما بر ۲ سطر اول A ، $u_1, u_2 \in \mathbb{R}^D$ و دو سطر اول در B ، $v_1, v_2 \in \mathbb{R}^k$ تمرکز می‌کنیم. تعریف می‌کنیم $R = \{r_{ij}\}_{i=1}^D \{j=1}^k$ بنابراین:

$$v_{1,j} = \sum_{i=1}^D r_{ij} u_{1,i}, \quad v_{2,j} = \sum_{i=1}^D r_{ij} u_{2,i}, \quad x_j = v_{1,j} - v_{2,j} = \sum_{i=1}^D r_{ij} (u_{1,i} - u_{2,i}). \quad (1-3)$$

۱-۱-۳ توزیع‌های پایدار

به طور معمول $r_{ij} \sim S(\alpha, 1)$ و به طور i.i.d. استخراج می‌شود. همچنین در ادامه ما حالت‌های ساده‌تری را هم مورد بررسی قرار می‌دهیم. در اینجا $S(\alpha, 1)$ بیانگر یک توزیع متقارن α -پایدار تصادفی است [۴۹] با پارامتر اندیس α و پارامتر مقیاس ۱. یک متغیر تصادفی z در صورتی متقارن α -پایدار است که تابع مشخصه‌ی آن به شکل زیر باشد.

$$E(\exp(\sqrt{-1}zt)) = \exp(-d|t|^\alpha) \quad (2-3)$$

که $d > 0$ پارامتر مقیاس است. ما می‌نویسیم $z \sim S(\alpha, d)$ که به طور کلی شکل بسته‌ای برای تابع چگالی ندارد. به جز حالت $\alpha = 2$ (نرمال) و $\alpha = 1$ (کوچی^۸).

^۶sub-Gaussian

^۷Cauchy random projection

^۸Cauchy

۳-۱-۲ مسئله برآورد آماری

با توجه به خواص تبدیل فوری، به راحتی می‌توان نشان داد که داده‌های تصویر شده هم از توزیع α -پایدار پیروی می‌کنند که در این حالت پارامتر مقیاس مشخصه l_α ی (نرم‌ها، فاصله‌ها) داده‌های اصلی در A است. به طور خاص:

$$v_{1,j} \sim S\left(\alpha, \sum_{i=1}^D |u_{1,i}|^\alpha\right), \quad v_{2,j} \sim S\left(\alpha, \sum_{i=1}^D |u_{2,i}|^\alpha\right), \quad (3-3)$$

$$x_j = v_{1,j} - v_{2,j} \sim S\left(\alpha, d_{(\alpha)} = \sum_{i=1}^D |u_{1,i} - u_{2,i}|^\alpha\right). \quad (4-3)$$

بنابراین، کار ما به برآورد پارامتر مقیاس از k نمونه i.i.d. $x_j \sim S(\alpha, d_{(\alpha)})$ ، تقلیل پیدا می‌کند. به این خاطر که هیچ شکل بسته‌ای برای تابع چگالی به جز در حالت $\alpha = 1, 2$ وجود ندارد، فرآیند تخمین خود مسئله‌ی جالبی است اگر به دنبال برآوردهایی بگردیم که هم به طور آماری دقیق باشند و هم از لحاظ محاسباتی کارا باشند.

یک موضوع مربوط و نزدیک هم تعیین اندازه نمونه k است. روش استاندارد محدود کردن احتمال دم است $\Pr(|\hat{d}_{(\alpha)} - d_{(\alpha)}| > \epsilon d_{(\alpha)})$ که $\hat{d}_{(\alpha)}$ برآوردگری برای $d_{(\alpha)}$ است و ϵ دقت مورد نظر است (معمولا $0 < \epsilon < 1$). به طور ایده‌آل امیدوار هستیم نشان دهیم ^۹:

$$\Pr(|\hat{d}_{(\alpha)} - d_{(\alpha)}| > \epsilon d_{(\alpha)}) \leq 2 \exp\left(-k \frac{\epsilon^2}{G}\right), \quad (5-3)$$

برای برخی مقادیر ثابت G که می‌تواند تابعی از ϵ هم باشد. برای ماتریس داده‌ی $A \in \mathbb{R}^{n \times D}$ ، در مجموع $\frac{n(n-1)}{2} < \frac{n^2}{2}$ جفت فاصله وجود دارد. ما معمولا علاقمندیم که احتمالات دم را به طور همزمان برای همه‌ی جفت‌ها محدود کنیم.

^۹ بنابر قضیه حد مرکزی برآوردگر $\hat{d}_{(\alpha)}$ بر اساس k نمونه تحت شروط ساده‌ای به حالت نرمال همگرا می‌شود. بنابر محدوده‌ی دم نرمال می‌دانیم که حداقل برای پارامترهای خاصی $\Pr(|\hat{d}_{(\alpha)} - d_{(\alpha)}| \geq \epsilon d_{(\alpha)}) \leq 2 \exp\left(-k \frac{\epsilon^2}{2V}\right)$ باید صادق باشد. در اینجا $\frac{V}{k}$ واریانس مجانبی $\hat{d}_{(\alpha)}$ است. بنابراین، حداقل برای آزمون درستی، می‌توانیم با بررسی این که آیا $\lim_{\epsilon \rightarrow 0+} G = 2V$ چک کنیم که محدوده‌ی دم نسبت مطلوب را دارا باشد.

۲-۳ تصویر تصادفی نرمال

برای کاهش بعد در نرم l_2 ، روش تصویر تصادفی نرمال ماتریس داده‌ی اولیه $A \in \mathbb{R}^{n \times D}$ را در ماتریس تصادفی $B \in \mathbb{R}^{n \times k}$ ($k \ll D$) با درایه‌های i.i.d. از $N(0, 1)$ ، تا ماتریس تصویر شده‌ی $B \in \mathbb{R}^{n \times k}$ حاصل شود. تحلیل‌های مربوط به تصویر تصادفی نرمال نسبتاً ساده است. برای مثال، به شکل سرراستی می‌توان یک نسخه از لم JL^{۱۰} [۳۵] را برای حالت l_2 استنتاج کرد.

ما در ابتدا برخی خواص اولیه تصویر تصادفی نرمال را بیان می‌کنیم و سپس بر روی اطلاعات حاشیه تمرکز می‌کنیم تا تخمین‌ها را بهینه کنیم. حاشیه‌ها (به عبارت دیگر، نرم l_2 برای هر خط در A) معمولاً در ابتدا در دسترس هستند (برای مثال، از طریق نرمال سازی داده‌ها). ولی حتی در حالتی که در دسترس نیستند، محاسبه‌ی نرم l_2 برای تمام سطرهای A فقط نیازمند یکبار مرور داده‌ها است که هزینه‌ای از $O(nD)$ دارد که قابل صرف‌نظر است.^{۱۱} از آنجا که اعمال تصویر تصادفی $A \times R$ هم اکنون هزینه‌ای از مرتبه $O(nDk)$ دارد.

در این بخش، ما این قاعده مرسوم تبعیت در ادبیات تصویر تصادفی [۴۷] پیروی می‌کنیم و تعریف می‌کنیم $B = \frac{1}{\sqrt{k}}AR$.

۱-۲-۳ مشخصه‌های اصلی

ما فرض می‌کنیم یک ماتریس داده $A \in \mathbb{R}^{n \times D}$ و یک ماتریس تصویرگر $R \in \mathbb{R}^{D \times k}$ که به طور i.i.d. از $N(0, 1)$ تولید شده است. در نظر می‌گیریم $B = \frac{1}{\sqrt{k}}AR$. در نظر بگیرید u_i^T سطر i ام ماتریس A باشد، و سطر متناظر در B ، v_i^T باشد. برای راحتی بر روی دو سطر اول A یعنی u_1 و u_2 همچنین دو سطر اولیه v_1 و v_2 در B تمرکز می‌کنیم. تعریف می‌کنیم:

$$a = u_1^T u_2, \quad m_1 = \|u_1\|^2, \quad m_2 = \|u_2\|^2, \quad d = \|u_1 - u_2\|^2 = m_1 + m_2 - 2a \quad (۶-۳)$$

به آسانی می‌توانیم نشان دهیم $\|v_1 - v_2\|$ ، فاصله‌ی l_2 نمونه و $v_1^T v_2$ ضرب داخلی نمونه، برآوردهای

¹⁰Johnson-Lindenstrauss

^{۱۱} این وضعیتی برای زمانی که با جریان داده‌های دینامیک سر و کار داریم اندکی متفاوت است. در جریان‌های داده ما معمولاً به دنبال اطلاعات آماری یک جریان داده هستیم تا اختلاف میان دو جریان داده را مد نظر داشته باشیم. به عبارت دیگر، محاسبه نرم l_2 حاشیه‌ای گاهی اوقات هدف اصلی است. به دلیل ذات دینامیک جریان‌های داده (برای مثال، به روز شدن مدام)، محاسبه‌ی حاشیه‌ها می‌تواند پر هزینه باشد.

نالریبی از d و a هستند. لم ۱ واریانس و تابع مشخصه‌ی $v_1^T v_2$ را مشخص می‌کند. اثبات در [۴۰].
 لم ۱: $u_1, u_2 \in \mathbb{R}^D$ داده شده‌اند و یک ماتریس تصادفی $\mathbf{R} \in \mathbb{R}^{D \times k}$ شامل درایه‌های i.i.d. از نرمال استاندارد $N(0, 1)$. اگر مقادیر $v_1 = \frac{1}{\sqrt{k}} \mathbf{R}^T u_1$ و $v_2 = \frac{1}{\sqrt{k}} \mathbf{R}^T u_2$ را تعیین کنیم، داریم:

$$E(\|v_1 - v_2\|^2) = d, \quad \text{Var}(\|v_1 - v_2\|^2) = \frac{2}{k} d^2 \quad (۷-۳)$$

$$E(v_1^T v_2) = a, \quad \text{Var}(v_1^T v_2) = \frac{1}{k} (m_1 m_2 + a^2), \quad (۸-۳)$$

سومین ممان مرکزی $v_1^T v_2$ عبارت است از:

$$E(v_1^T v_2)^2 = a, \quad \frac{2a}{k^2} (2m_1 m_2 + a^2) \quad (۹-۳)$$

و تابع مولد احتمال برای $v_1^T v_2$ عبارت است از:

$$E(\exp(v_1^T v_2 t)) = \left(1 - \frac{2}{k} a t - \frac{1}{k^2} (m_1 m_2 - a^2) t^2\right)^{-\frac{k}{2}} \quad (۱۰-۳)$$

که $\frac{-k}{\sqrt{m_1 m_2 - a}} \leq t \leq \frac{-k}{\sqrt{m_1 m_2 + a}}$ است.

بنابراین، برآوردهای نالریبی برای فاصله d و ضرب داخلی a به شکل سر راستی عبارت است از:

$$\hat{d}_{MF} = \|v_1 - v_2\|^2, \quad \text{Var}(\hat{d}_{MF}) = \frac{d^2}{k}, \quad (۱۱-۳)$$

$$\hat{a}_{MF} = v_1^T v_2, \quad \text{Var}(\hat{a}_{MF}) = \frac{1}{k} (m_1 m_2 + a^2), \quad (۱۲-۳)$$

که اندیس « MF » به معنی «بدون حاشیه»^{۱۲} نشان دهنده این است که برآوردها از اطلاعات حاشیه‌ای $m_1 = \|u_1\|^2$ و $m_2 = \|u_2\|^2$ استفاده نمی‌کنند.

به این نکته توجه کنید که، $k \hat{d}_{MF} / d$ از توزیع χ^2 با k درجه آزادی، پیروی می‌کند، χ_k^2 . بنابراین،

^{۱۲}margin-free

به راحتی می‌توان می‌توانیم این محدوده‌های دم را برای لم ۲ اثبات کنیم.

لم ۲:

$$\Pr(\hat{d}_{MF} - d > \epsilon d) \leq \exp\left(-\frac{k}{2}(\epsilon - \log(1 + \epsilon))\right), \quad \epsilon > 0 \quad (۱۳-۳)$$

$$\Pr(\hat{d}_{MF} - d < -\epsilon d) \leq \exp\left(-\frac{k}{2}(-\epsilon - \log(1 - \epsilon))\right), \quad 0 < \epsilon < 1 \quad (۱۴-۳)$$

اثبات:

از آنجا که $k\hat{d}_{MF}/d \sim \chi_k^2$ ، بر اساس نام مساوی چرنوف^{۱۳} [۱۹]، برای هر $t > 0$ داریم:

$$\begin{aligned} \Pr(\hat{d}_{MF} - d > \epsilon d) &= \Pr(k\hat{d}_{MF}/d > k(1 + \epsilon)) \\ &\leq \frac{E\left(\exp(k\hat{d}_{MF}/dt)\right)}{\exp((1 + \epsilon)kt)} = \exp\left(-\frac{k}{2}(\log(1 - 2t) + 2(1 + \epsilon)t)\right) \end{aligned} \quad (۱۵-۳)$$

که در $t = t_{NR} = \frac{\epsilon}{2(1+\epsilon)}$ و بنابراین برای هر $\epsilon > 0$ داریم:

$$\Pr(\hat{d}_{MF} - d > \epsilon d) \leq \exp\left(-\frac{k}{2}(\epsilon - \log(1 + \epsilon))\right) \quad (۱۶-۳)$$

ما می‌توانیم به طور مشابه برای دیگر محدوده‌ی دم $\Pr(\hat{d}_{MF} - d < -\epsilon d)$ هم اثبات کنیم. ■

برای راحتی مرسوم است که محدوده دم را در لم ۲ به صورت متقارن $\Pr(|\hat{d}_{MF} - d| > \epsilon d)$ نوشته شود. نامساوی‌های ساده‌ای برای $\log(1 + \epsilon)$ و $\log(1 - \epsilon)$ نتیجه می‌دهد:

$$\Pr(|\hat{d}_{MF} - d| \geq \epsilon d) \leq 2 \exp\left(-\frac{k}{4}\epsilon^2 + \frac{k}{6}\epsilon^3\right), \quad 0 < \epsilon < 1 \quad (۱۷-۳)$$

از آنجا که $\mathbf{A} \in \mathbb{R}^{n \times D}$ تعداد n سطر دارد. به عبارت دیگر $\frac{n(n-1)}{2}$ جفت. ما باید احتمال دم را به

¹³Chernoff inequality

طور همزمان برای همه‌ی جفت‌ها محدود کنیم. با استفاده از محدوده تجمیعی بنفرونی^{۱۴} کافی است که:

$$\frac{n^2}{2} \Pr \left(\left| \hat{d}_{MF} - d \right| \geq \epsilon d \right) \leq \delta \quad (۱۸-۳)$$

به عبارت دیگر کافی است اگر:

$$\frac{n^2}{2} 2 \exp \left(-\frac{k}{4} \epsilon^2 + \frac{k}{6} \epsilon^3 \right) \leq \delta \Rightarrow k \geq \frac{2 \log n - \log \delta}{\epsilon^2/4 - \epsilon^3/6} \quad (۱۹-۳)$$

بنابراین ما یک نسخه‌ای از لم JL را نشان داده‌ایم.

لم ۳: اگر $k \geq \frac{2 \log n - \log \delta}{\epsilon^2/4 - \epsilon^3/6}$ پس با حداقل احتمال $1 - \delta$ ، فاصله l_2 بین هر جفت از داده‌ها (میان n نقطه) می‌تواند با ضریب اطمینان $1 \pm \epsilon$ با استفاده فاصله‌ی l_2 در داده‌های تصویر شده بعد از تصویر تصافی نرمال، تخمین زده شود. $0 < \delta < 1, 0 < \epsilon < 1$. ■

۳-۳ تصویر تصادفی زیر گوسی و بسیار پراکنده

در بخش قبل ما به بررسی تصویر تصادفی نرمال پرداختیم، که در آن ماتریس تصویر گر R از روی توزیع $N(0, 1)$ به طور i.i.d. نمونه‌گیری می‌شود. این انتخاب خاص برای R ، صرفاً برای سهولت تحلیل تئوری است. در واقع می‌توان R را از هر توزیعی با میانگین صفر و واریانس محدود برای کاهش بعد در نرم l_2 نمونه‌گیری کرد.

نمونه‌گیری R از یک توزیع زیرگوسی^{۱۵} هم از نظر تئوری قابل قبول و هم از جنبه‌ی محاسباتی تسهیل کننده است. برای مثال، محدوده‌ی دم زیر گوسی به سادگی به نسخه‌ای از لم JL منتهی می‌شود. ما بر روی یک انتخاب معمول از توزیع زیر گوسی تمرکز خواهیم کرد، که درایه‌ها ماتریس R از مجموعه‌ی $\{-1, 0, 1\}$ با احتمالات $\{\frac{1}{2s}, 1 - \frac{1}{s}, \frac{1}{2s}\}$ ، که $s \geq 1$. به این ترتیب فرآیند نمونه‌گیری ساده‌تر شده و محاسبات سریعتر انجام می‌شوند. در واقع، زمانی که $s < 3$ باشد، واریانس‌های صرحا

^{۱۴}Benferroni union bound

^{۱۵}sub-Gaussian

کوچکتری نسبت به استفاده از تصویر تصادفی نرمال بدست می‌آید.
 با در نظر گرفتن قواعد معقول، برای مثال، داده‌های اولیه ممان سوم محدود داشته باشند، می‌توانیم
 $s \gg 3$ در نظر بگیریم (حتی $s = \sqrt{D}$). تا نتایج s برابر سریعتر بدست بیاوریم؛ و بنابراین، این رویه را
 تصویر تصادفی بسیار پراکنده می‌نامیم.

۳-۳-۱ تصویر تصادفی زیرگوسی

مشابه **قسمت ۲-۳** ماتریس داده را $A \in \mathbb{R}^{n \times D}$ در نظر می‌گیریم. ماتریس تصویر تصادفی $R \in \mathbb{R}^{D \times k}$
 را تولید کرده و آن را در A ضرب می‌کنیم تا به یک ماتریس تصویر شده $B = \frac{1}{\sqrt{k}} AR \in \mathbb{R}^{n \times k}$
 برسیم. دوباره رو دو ردیف ابتدایی تمرکز می‌کنیم، که یعنی u_1 و u_2 در A ، و دو ردیف ابتدایی v_1 و v_2
 در B و همچنین تساوی‌های زیر را تعریف می‌کنیم:

$$a = u_1^T u_2, \quad m_1 = \|u_1\|^2, \quad m_2 = \|u_2\|^2, \quad d = \|u_1 - u_2\|^2 = m_1 + m_2 - 2a \quad (۲۰-۳)$$

R را به طور i.i.d از یک توزیع زیر گوسی مشخصا پرکاربرد تولید می‌کنیم: ($S \geq 1$)

$$r_{ij} = \sqrt{s} \times \begin{cases} 1 & \text{با احتمال } \frac{1}{2s} \\ 0 & \text{با احتمال } 1 - \frac{1}{s} \\ -1 & \text{با احتمال } \frac{1}{2s} \end{cases} \quad (۲۱-۳)$$

- نمونه‌گیری از **معادله ۲۱-۳** ساده‌تر از نمونه‌گیری از $N(0, 1)$ است.
- می‌تواند از s برابر افزایش سرعت در ضرب ماتریسی $A \times R$ بهره برد، زیرا فقط $\frac{1}{s}$ داده‌های نیازمند پردازش هستند.
- نیازی به عملیات محاسباتی با ممیز شناور نیست و تمامی بار محاسباتی بر روی عملیات تجمیع پایگاه داده است که به خوبی بهینه شده.
- وقتی $s < 3$ باشد می‌توان به تخمین‌هایی با دقت بیشتر (واریانس کمتر) دست پیدا کرد.

- هزینه نگهداری ماتریس \mathbf{R} از $O(Dk)$ به $O(Dk/s)$ کاهش می‌یابد.

[۱، ۲] نشان می‌دهند زمانی که $s = 1$ و $s = 3$ باشد، می‌توان به همان محدوده‌ی JL ای دست پیدا کرد که در تصویر تصادفی نرمال وجود دارد. ما در ادامه به بررسی خواص توزیع زیر گوسی می‌پردازیم، که برای تحلیل محدوده‌ی دم مناسب است. در واقع، آنالیز زیر گوسی نشان می‌دهد که می‌توان حتی در بدترین شرایط از مقادیری اندکی بیشتر از ۳ برای s استفاده کرد.

توزیع زیر گوسی

ما در اینجا مقدمه‌ای کوتاه بر توزیع‌های زیرگوسی بیان می‌کنیم. برای جزئیات و منابع بیشتر می‌توانید به [۱۷] مراجعه کنید. تئوری توزیع‌های زیرگوسی در حدود ۱۹۶۰ آغاز شد. متغیر تصادفی x زیرگوسی است اگر ثابت $g > 0$ وجود داشته باشد به شکلی که:

$$\mathbb{E}(\exp(xt)) \leq \exp\left(\frac{g^2 t^2}{2}\right), \forall t \in \mathbb{R} \quad (۲۲-۳)$$

می‌توان مقدار بهینه‌ی g^2 را از تعریف $T^2(x)$ با استفاده از فرمول زیر بدست آورد.

$$T^2(x) = \sup_{t \neq 0} \frac{2 \log \mathbb{E}(\exp(xt))}{t^2} \quad (۲۳-۳)$$

توجه کنید که $T^2(x)$ فقط یک نمادگذاری برای مقدار ثابت بهینه‌ی زیرگوسی یک متغیر تصادفی x است (و نه یک نمونه مشخص از x).

برخی از ویژگی‌های اولیه‌ی توزیع‌های زیرگوسی:

- اگر x زیر گوسی باشد آنگاه $\mathbb{E}(x) = 0$ و $\mathbb{E}(x^2) \leq T^2(x)$. برای هر مقدار ثابت c ، $T^2(cx) = c^2 T^2(x)$ و

$$\Pr(|x| > t) \leq 2 \exp\left(-\frac{t^2}{2T^2(x)}\right) \quad (۲۴-۳)$$

- اگر x_1, x_2, \dots, x_D زیرگوسی مستقل باشند، آنگاه $\sum_{i=1}^D x_i$ زیرگوسی است.

$$T^2 \left(\sum_{i=1}^D x_i \right) \leq \sum_{i=1}^D T^2(x_i) \quad (25-3)$$

• اگر x زیرگوسی باشد، آنگاه برای همه $t \in [0, 1)$ ،

$$\mathbf{E} \left(\exp \left(\frac{x^2 t}{2T^2(x)} \right) \right) \leq (1 - t)^{-\frac{1}{2}} \quad (26-3)$$

[۱، ۲] همچنین معادله ۲۶-۳ را برای توزیع ویژه‌ی معادله ۲۱-۳ بدست آورده‌اند. یک متغیر

تصادفی زیرگوسی x صریحا زیرگوسی است اگر $\mathbf{E}(x^2) = T^2(x)$

• اگر x صریحا زیرگوسی باشد، آنگاه $\mathbf{E}(x^3) = 0$ و کشیدگی^{۱۶} غیر مثبت خواهد بود، به عبارت

$$\text{دیگر } \frac{\mathbf{E}(x^4)}{\mathbf{E}^2(x^2)} - 3 \leq 0$$

• اگر x_1, x_2, \dots, x_D صریحا زیرگوسی مستقل باشند، آنگاه $\sum_{i=1}^D x_i$ صریحا زیرگوسی است.

$$T^2 \left(\sum_{i=1}^D x_i \right) = \sum_{i=1}^D T^2(x_i) = \sum_{i=1}^D \mathbf{E}(x_i^2) \quad (27-3)$$

۴-۳ تصویر تصادفی کوچی

توضیحات

۵-۳ تصویر تصادفی α -پایدار

توضیحات

¹⁶kurtosis

فصل چهارم

کاهش بعد و نحوه‌ی بررسی عملکرد آن

توضیحات

۱-۴ PCA و مقایسه با آن

توضیحات

۲-۴ برآوردگرای معیارهای وابستگی

توضیحات

۳-۴ داده‌های مورد استفاده

توضیحات

۴-۴ توضیحات کد

توضیحات

۵-۴ شاخص محاسبه عملکرد کاهش بعد و Adjusted Rand

Index

توضیحات

فصل پنجم

نتایج

منابع و مراجع

- [1] Achlioptas, Dimitris. Database-friendly random projections. In *Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 274–281. ACM, 2001.
- [2] Achlioptas, Dimitris. Database-friendly random projections: Johnson-lindenstrauss with binary coins. *Journal of computer and System Sciences*, 66(4):671–687, 2003.
- [3] Aggarwal, Charu C. *Data streams: models and algorithms*, volume 31. Springer Science & Business Media, 2007.
- [4] Aggarwal, Charu C, Wolf, Joel L, and Yu, Philip S. *A new method for similarity indexing of market basket data*. ACM, 1999.
- [5] Agrawal, Rakesh, Imieliński, Tomasz, and Swami, Arun. Mining association rules between sets of items in large databases. In *Acm sigmod record*, volume 22, pages 207–216. ACM, 1993.
- [6] Agrawal, Rakesh, Mannila, Heikki, Srikant, Ramakrishnan, Toivonen, Hannu, Verkamo, A Inkeri, et al. Fast discovery of association rules. *Advances in knowledge discovery and data mining*, 12(1):307–328, 1996.
- [7] Agrawal, Rakesh, Srikant, Ramakrishnan, et al. Fast algorithms for mining association rules. In *Proc. 20th int. conf. very large data bases, VLDB*, volume 1215, pages 487–499, 1994.
- [8] Andoni, Alexandr and Indyk, Piotr. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pages 459–468. IEEE, 2006.

- [9] Babcock, Brian, Babu, Shivnath, Datar, Mayur, Motwani, Rajeev, and Widom, Jennifer. Models and issues in data stream systems. In *Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 1–16. ACM, 2002.
- [10] Brin, Sergey, Davis, James, and Garcia-Molina, Hector. Copy detection mechanisms for digital documents. In *ACM SIGMOD Record*, volume 24, pages 398–409. ACM, 1995.
- [11] Brin, Sergey, Motwani, Rajeev, and Silverstein, Craig. Beyond market baskets: Generalizing association rules to correlations. In *Acm Sigmod Record*, volume 26, pages 265–276. ACM, 1997.
- [12] Brin, Sergey, Motwani, Rajeev, Ullman, Jeffrey D, and Tsur, Shalom. Dynamic itemset counting and implication rules for market basket data. *Acm Sigmod Record*, 26(2):255–264, 1997.
- [13] Brin, Sergey and Page, Lawrence. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1-7):107–117, 1998.
- [14] Brinkman, Bo and Charikar, Moses. On the impossibility of dimension reduction in ℓ_1 . *Journal of the ACM (JACM)*, 52(5):766–788, 2005.
- [15] Broder, Andrei Z. On the resemblance and containment of documents. In *Compression and complexity of sequences 1997. proceedings*, pages 21–29. IEEE, 1997.
- [16] Buhler, Jeremy and Tompa, Martin. Finding motifs using random projections. *Journal of computational biology*, 9(2):225–242, 2002.
- [17] Buldygin, Valeri Vladimirovich and Kozachenko, IU V. *Metric characterization of random variables and random processes*, volume 188. American Mathematical Soc., 2000.
- [18] Chaudhuri, Surajit, Motwani, Rajeev, and Narasayya, Vivek. On random sampling over joins. In *ACM SIGMOD Record*, volume 28, pages 263–274. ACM, 1999.

- [19] Chernoff, Herman et al. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, 23(4):493–507, 1952.
- [20] Crovella, Mark E and Bestavros, Azer. Self-similarity in world wide web traffic: evidence and possible causes. *IEEE/ACM Transactions on networking*, 5(6):835–846, 1997.
- [21] Datar, Mayur, Immorlica, Nicole, Indyk, Piotr, and Mirrokni, Vahab S. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on Computational geometry*, pages 253–262. ACM, 2004.
- [22] Datar, Mayur and Indyk, Piotr. Comparing data streams using hamming norms. In *Proceedings 2002 VLDB Conference: 28th International Conference on Very Large Databases (VLDB)*, page 335. Elsevier, 2002.
- [23] Deerwester, Scott, Dumais, Susan T, Furnas, George W, Landauer, Thomas K, and Harshman, Richard. Indexing by latent semantic analysis. *Journal of the American society for information science*, 41(6):391–407, 1990.
- [24] Dhillon, Inderjit S and Modha, Dharmendra S. Concept decompositions for large sparse text data using clustering. *Machine learning*, 42(1-2):143–175, 2001.
- [25] Faloutsos, Michalis, Faloutsos, Petros, and Faloutsos, Christos. On power-law relationships of the internet topology. In *ACM SIGCOMM computer communication review*, pages 251–262. ACM, 1999.
- [26] Friedman, Jerome, Hastie, Trevor, and Tibshirani, Robert. *The elements of statistical learning*, volume 10. Springer series in statistics New York, NY, USA:, 2001.
- [27] Friedman, Jerome H, Baskett, Forest, and Shustek, Leonard J. An algorithm for finding nearest neighbors. *IEEE Transactions on computers*, 100(10):1000–1006, 1975.
- [28] Friedman, Jerome H, Bentley, Jon Louis, and Finkel, Raphael Ari. An algorithm for finding best matches in logarithmic time. *ACM Trans. Math. Software*, 3(SLAC-PUB-1549-REV. 2):209–226, 1976.

- [29] Garcia-Molina, Hector. Database systems: the complete book/hector garcia, molina jeffrey d. ullman, jennifer widom, 2002.
- [30] Henzinger, Monika Rauch, Raghavan, Prabhakar, and Rajagopalan, Sridhar. Computing on data streams. *External memory algorithms*, 50:107–118, 1998.
- [31] Hornby, Albert Sydney, editor. *Oxford Advanced Learner's Dictionary of Current English*. Oxford University Press, Oxford, UK, fourth edition, 1989.
- [32] Indyk, Piotr. Stable distributions, pseudorandom generators, embeddings and data stream computation. In *focs*, page 189. IEEE, 2000.
- [33] Indyk, Piotr. Stable distributions, pseudorandom generators, embeddings, and data stream computation. *Journal of the ACM (JACM)*, 53(3):307–323, 2006.
- [34] Indyk, Piotr and Motwani, Rajeev. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613. ACM, 1998.
- [35] Johnson, William B and Lindenstrauss, Joram. Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics*, 26(189-206):1, 1984.
- [36] Johnson, William B and Schechtman, Gideon. Embedding l_p into l_1 . *Acta Mathematica*, 149(1):71–85, 1982.
- [37] Kannan, J Feigenbaum S, Strauss, M, and Viswanathan, M. An approximate l_1 -difference algorithm for massive data streams. *Unknown*, Unknown.
- [38] Lee, James R and Naor, Assaf. Embedding the diamond graph in l_p and dimension reduction in l_1 . *Geometric & Functional Analysis GAFA*, 14(4):745–747, 2004.
- [39] Leland, Will E, Willinger, Walter, Taqqu, Murad S, and Wilson, Daniel V. On the self-similar nature of ethernet traffic. *ACM SIGCOMM Computer Communication Review*, 25(1):202–213, 1995.
- [40] Li, Ping. *Stable random projections and conditional random sampling, two sampling techniques for modern massive datasets*. Stanford, 2007.

- [41] Li, Ping. Estimators and tail bounds for dimension reduction in l_α ($0 < \alpha \leq 2$) using stable random projections. In *Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 10–19. Society for Industrial and Applied Mathematics, 2008.
- [42] Matias, Yossi, Vitter, Jeffrey Scott, and Wang, Min. Wavelet-based histograms for selectivity estimation. In *ACM SIGMOD Record*, volume 27, pages 448–459. ACM, 1998.
- [43] McKee, Sally A. Reflections on the memory wall. In *CF'04: Proceedings of the 1st conference on Computing frontiers*, page 162, 2004.
- [44] Muthukrishnan, S. Data streams: Algorithms and applications (foundations and trends in theoretical computer science). *Hanover, MA: Now Publishers Inc*, 2005.
- [45] Newman, Mark EJ. Power laws, pareto distributions and zipf's law. *Contemporary physics*, 46(5):323–351, 2005.
- [46] Strehl, Alexander and Ghosh, Joydeep. A scalable approach to balanced, high-dimensional clustering of market-baskets. In *International Conference on High-Performance Computing*, pages 525–536. Springer, 2000.
- [47] Vempala, Santosh S. *The random projection method*, volume 65. American Mathematical Soc., 2005.
- [48] Wulf, Wm A and McKee, Sally A. Hitting the memory wall: implications of the obvious. *ACM SIGARCH computer architecture news*, 23(1):20–24, 1995.
- [49] Zolotarev, VM. One-dimensional stable distributions. translated from the russian by hh mcfaden. translation edited by ben silver. translations of mathematical monographs, 65. *American Mathematical Society, Providence, RI*, 1986.

پیوست

موضوعات مرتبط با متن گزارش پایان نامه که در یکی از گروه‌های زیر قرار می‌گیرد، در بخش پیوست‌ها آورده شوند:

۱. اثبات‌های ریاضی یا عملیات ریاضی طولانی.

۲. داده و اطلاعات نمونه (های) مورد مطالعه (Case Study) چنانچه طولانی باشد.

۳. نتایج کارهای دیگران چنانچه نیاز به تفصیل باشد.

۴. مجموعه تعاریف متغیرها و پارامترها، چنانچه طولانی بوده و در متن به انجام نرسیده باشد.

کد میپل

```
with(DifferentialGeometry):  
with(Tensor):  
DGsetup([x, y, z], M)  
frame name: M  
a := evalDG(D_x)  
D_x  
b := evalDG(-2 y z D_x+2 x D_y/z^3-D_z/z^2)
```

واژه‌نامه‌ی فارسی به انگلیسی

آ	Cartesian product حاصل ضرب دکارتی
اسکالر Scalar	خ
ب	Automorphism خودریختی
بالابر Lift	د
پ	Degree درجه
پایا Invariant	ر
ت	microprocessor ریزپردازنده
تناظر Correspondence	ز
ث	Submodule زیرمدول
ثابت‌ساز Stabilizer	س
ج	Character سرشت
جایگشت Permutation	ص
چ	Faithful صادقانه
چند جمله‌ای Polynomial	ض
ح	

Connected	همبند	Inner product	ضرب داخلی
	ی		ط
Edge	یال	Loop	طوقه
			ظ
		Valency	ظرفیت
			ع
		Nonadjacency	عدم مجاورت
			ف
		Vector space	فضای برداری
			ک
		Complete reducibility	کاملاً تحویل پذیر
			گ
		Graph	گراف
			م
		Permutation matrix	ماتریس جایگشتی
			ن
		Disconnected	ناهمبند
			و
		Invertible	وارون پذیر
			ه

واژه‌نامه‌ی انگلیسی به فارسی

A	Homomorphism همریختی
Automorphism خودریختی	I
B	Invariant پایا
Bijection دوسویی	L
C	Lift بالابر
Cycle group گروه دوری	M
D	Module مدول
Degree درجه	N
E	Natural map نگاشت طبیعی
Edge یال	O
F	One to One یک به یک
Function تابع	P
G	Permutation group گروه جایگشتی
Group گروه	Q
H	

Quotient graph	گراف خارج‌قسمتی	Trivial character	سرشت بدیهی
R		U	
Reducible	تحویل پذیر	Unique	منحصربفرد
S		V	
Sequence	دنباله	Vector space	فضای برداری
T			

Abstract

This page is accurate translation from Persian abstract into English.

Key Words:

Write a 3 to 5 KeyWords is essential. Example: AUT, M.Sc., Ph. D, ..



Amirkabir University of Technology
(Tehran Polytechnic)

Department of ...

M. Sc. Thesis

Title of Thesis

By

Name Surname

Supervisor

Dr.

Advisor

Dr.

Month & Year