



الگوریتم‌های خلاصه‌سازی برای مه‌داده

محمد هادی فروغ‌منداغرابی
پاییز ۱۳۹۹

خلاصه‌سازی خطی برای تقریبی از آرایه با افزایش و کاهش خانه‌ها

جلسه هفتم

نگارنده: امیر آذر مهر

۱ صورت مسئله

فرض کنید یک بردار $x \in \mathbb{R}^n$ داریم، که طی زمان تغییراتی می‌کند. می‌خواهیم با استفاده از حافظه‌ی $n \ll m$ ، تقریبی از مقدار x را ذخیره کنیم. به صورت دقیق‌تر برای نوع تغییرات x سه حالت در نظر گرفته می‌شود:

- فقط اضافه کردن؛ یعنی با هر تغییر به ازای یک i قرار می‌دهیم $x_i \leftarrow x_i + 1$
 - شمارنده‌ی مثبت^۱؛ یعنی هر بار قرار می‌دهیم $x_i \leftarrow x_i \pm 1$ ، به شکلی که تمام مؤلفه‌ها همواره نامنفی باشند.
 - شمارنده‌ی آزاد^۲؛ یعنی هر بار قرار می‌دهیم $x_i \leftarrow x_i \pm 1$ و محدودیتی روی مقدار مؤلفه‌ها وجود ندارد.
- همچنین برای نگهداری تقریبی x یکی از دو شیوه‌ی زیر را برای پاسخگویی الگوریتم در نظر می‌گیریم:
- پرسش نقطه‌ای (مؤلفه‌ای)^۳ ℓ_1 : در این شیوه پارامتر k یک عدد ثابت است. اندیس i به الگوریتم داده می‌شود و انتظار می‌رود که مقدار x_i را با حداکثر خطای مطلق $\|x\|_1/k$ گزارش دهد.
 - پرسش وزین‌ها^۴ ℓ_1 : در این شیوه هم پارامتر k ثابت است و از الگوریتم انتظار می‌رود که هنگام پرسش، یک زیرمجموعه‌ی S از اندیس‌ها خروجی دهد. به شکلی که اندازه‌ی S از مرتبه $O(k)$ باشد و تمام اندیس‌هایی مثل i که $|x_i| \geq \|x\|/k$ در S آمده باشند. نکته: دقت کنید

^۱strict turnstile
^۲general turnstile
^۳point query
^۴heavy hitters

مجموعه‌ی اندیس‌های دارای خاصیت گفته شده، همواره حداکثر k عضو دارد و الگوریتم از این جهت تقریبی است که اجازه دارد $O(k)$ اندیس بازگرداند. این آزادی بیشتر اجازه می‌دهد که حافظه‌ی کمتری مصرف شود.

۲ کاهش پرسش وزین‌ها به پرسش نقطه‌ای

قضیه‌ای که در ادامه می‌آید، به مفهومی نشان می‌دهد که مسئله پرسش نقطه‌ای سخت‌تر از پرسش وزین است. توجه داشته باشید که الگوریتم‌هایی که مورد بررسی هستند همگی تصادفی هستند و به احتمال مشخصی پرسش را درست جواب می‌دهند.

قضیه ۱. اگر الگوریتم A برای پاسخ دادن پرسش نقطه‌ای ℓ_1 با پارامتر $3k$ وجود داشته باشد که هر پرسش را به احتمال حداکثر کمتر از δ/n اشتباه پاسخ دهد. آنگاه الگوریتمی برای پاسخ دادن پرسش وزین‌ها ℓ_1 با پارامتر k وجود دارد که هر پرسش را به احتمال کمتر از δ اشتباه پاسخ دهد.

اثبات. برای ساخت الگوریتم A' که پرسش وزین‌ها را پاسخ دهد به این شکل عمل می‌کنیم: از الگوریتم A برای اعمال تغییرها استفاده می‌کنیم. برای پاسخ دادن پرسش، مجموعه‌ی S را برابر با $3k$ بزرگترین مؤلفه، طبق تقریبی که A از آن‌ها دارد، قرار می‌دهیم. باید ثابت کنیم که S ، به احتمال مورد نظر، دو خاصیت مورد نظر را دارد.

این پدیده که A پرسش i -ام را اشتباه پاسخ دهد را E_i می‌نامیم. طبق فرض احتمال وقوع E_i حداکثر δ/n است. پس طبق کران اجتماع احتمال این که حداقل یکی از پرسش‌ها اشتباه پاسخ داده شود کمتر از δ است. پس اینکه همه‌ی پرسش‌ها به درستی توسط A پاسخ داده شوند، حداقل $1 - \delta$ است. نشان می‌دهیم در این صورت S خواص مورد نظر را دارد.

اندازه‌ی S طبق تعریف از مرتبه‌ی $O(k)$ است. برای اثبات خاصیت دوم، مقدار $t = \frac{2}{3k} \|x\|$ را در نظر بگیرید. اثبات می‌کنیم تمام مؤلفه‌هایی که مقدار گزارش شده‌ی آن‌ها (منظور تقریبی است که A خروجی می‌دهد) از t بزرگتر مساوی باشد در S قرار خواهد گرفت. چون که اگر به ازای i داشته باشیم $|x_i| < \frac{1}{3k} \|x\|$ آنگاه مقدار گزارش شده‌ی آن، طبق خواص الگوریتم A ، مقدار کوچکتر از $t = \frac{2}{3k} \|x\|$ دارد. پس از آنجا که حداکثر $3k$ اندیس i موجود هستند به طوری که $|x_i| \geq \frac{1}{3k} \|x\|$ ، تعداد اندیس‌ها با مقدار گزارش شده‌ی بزرگتر مساوی t هم حداکثر $3k$ است. پس تمامی این اعداد در S ظاهر می‌شوند. حالا در نظر داشته باشید که به ازای هر i که $|x_i| \geq \frac{1}{k} \|x\|$ ، باز هم طبق خواص الگوریتم A ، مقدار گزارش شده حداقل t است. پس تمام این اعداد در S قرار می‌گیرند. که خاصیت دوم مجموعه و در نتیجه درستی الگوریتم را به اثبات می‌رساند. \square

۳ الگوریتم COUNTMIN: پرسش نقطه‌ای برای شمارنده‌ی مثبت

به صورت کلی یک الگوریتم ارائه می‌دهیم که با احتمال ثابت پاسخ درست می‌دهد. با ترکیب خروجی اجراهای موازی این الگوریتم، یک پاسخ بدست می‌آوریم که به احتمال دلخواه درست است. ماتریس $C \in \mathbb{R}^{L \times B}$ که در ابتدا تمام صفر است در نظر بگیرید (هر سطر مستقل از بقیه سطرها است و قرار است به احتمال ثابت پاسخ درست دهد). در اینجا $B = 2k$ و $L = \lceil \log(1/\delta) \rceil$ همچنین برای هر سطر a یک تابع درهم‌سازی 2 -طرفه h_a از $[n]$ به $[B]$ در نظر بگیرید، به طوری که تابع‌های سطرها متفاوت از هم کاملاً مستقل باشند. حالا به ازای تغییر Δ $x_i \leftarrow x_i + \Delta$ برای هر سطر a مؤلفه‌ی $h_a(i)$ -ام را با Δ جمع می‌کنیم. در نهایت وقتی مقدار x_i درخواست می‌شود مقدار $C_{ah_a(i)}$ را گزارش می‌دهیم.

۱.۳ اثبات درستی

برای گزارش مقدار x_i ابتدا مقدار $C_{ah_a(i)}$ را برای یک سطر بررسی می‌کنیم. این مقدار برابر است با جمع تمام مؤلفه‌هایی از x که توسط تابع h_a به خانه‌ی $C_{ah_a(i)}$ وصل می‌شوند. به این ترتیب اگر Z_j متغیر تصادفی مشخصه‌ی این پدیده در نظر بگیرید که $h_a(j) = h_a(i)$ آنگاه داریم:

$$C_{ah_a(i)} = x_i + \underbrace{\sum_{j \neq i} Z_j x_j}_E$$

در اینجا E مقدار خطای محاسبه‌ی x_i و همواره یک عدد نامنفی است. امید ریاضی آن را حساب می‌کنیم:

$$\begin{aligned} \mathbb{E}[E] &= \sum_{j \neq i} \mathbb{E}[Z_j] x_j \\ &\leq \frac{1}{B} \sum_{j \neq i} x_j \\ &\leq \frac{1}{B} \|x\|_1 \end{aligned}$$

union bound^۵

که در اینجا خط دوم از ۲- طرفه بودن تابع درهم‌سازی نتیجه می‌شود. حالا با استفاده از نامساوی مارکوف یک کران بالا برای احتمال اینکه خطا بیش از حد مجاز باشد ارائه می‌دهیم:

$$\mathbb{P}[E > \|x_j\|/k] \leq \frac{\mathbb{E}[E]}{\|x_j\|/k} \leq \frac{k}{B} = \frac{1}{2}$$

به این ترتیب هر $Cha(i)$ به ازای هر سطر a با احتمال $\frac{1}{2}$ یک تقریب قابل قبول است. از آنجا که خطاها همیشه مثبت هستند، مقدار $\min_a Cha(i)$ در صورتی غیر قابل قبول است که همه‌ی سطرها غیر قابل قبول باشند. این اتفاق به احتمال کمتر از $\delta = 2^{-L}$ رخ می‌دهد. که درستی الگوریتم را اثبات می‌کند.

۲.۳ تحلیل زمان اجرا و حافظه

برای هر تغییر یا هر پرسش نقطه‌ای، باید عدد ورودی را بخوانیم و یک مؤلفه از هر سطر را تغییر دهیم یا بررسی کنیم. زمان مربوط به این کار $\mathcal{O}(\log(n/\delta)) + \mathcal{O}(\log(1/\delta)) = \mathcal{O}(\log(n/\delta))$ است. همچنین باید کل ماتریس را همواره نگهداری کنیم که حافظه‌ی مصرفی از مرتبه‌ی $\mathcal{O}(k \log(1/\delta))$ دارد. موقع پاسخ دادن پرسش‌ها باید k اندیس را نگهداری کنیم که حافظه‌ی $\mathcal{O}(k \log(n))$ نیاز دارد. پس در کل حافظه‌ی مصرفی $\mathcal{O}(k \log(n/\delta))$ است.

نکته: توجه کنید حافظه‌ی لازم برای نگهداری توابع درهم‌سازی بررسی نشده است و در بخش‌های بعد هم بررسی نمی‌شود.

۳.۳ پرسش وزین‌ها

طبق قضیه ۲؟ با همین حافظه، زمان $\mathcal{O}(\log(n/\delta))$ برای تغییر و زمان $\mathcal{O}(n \log(n/\delta))$ برای پرسش‌ها، می‌توان پرسش وزین‌ها را حل کرد. زمان اجرای مربوط به پرسش‌ها خیلی زیاد است که در بخش بعد به آن می‌پردازیم.

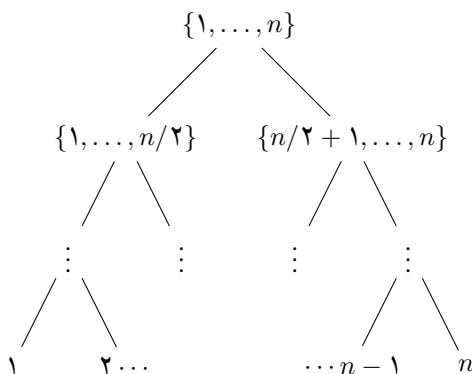
۴.۳ دلیل نام‌گذاری فشرده‌سازی خطی

دلیل این نام‌گذاری به این خاطر است که هر عدد که در ماتریس نگه می‌داریم، یک ترکیب خطی از اعداد اصلی است. به صورت کلی در این روش‌ها یک ماتریس تصادفی A داریم و به جای ذخیره کردن x مقدار $y = Ax$ را نگهداری می‌کنیم. برای تغییر $x \leftarrow x + e_i$ کافی است قرار درهم‌سازی $y \leftarrow y + Ae_i$. پاسخ پرسش‌ها به شیوه‌ای از روی y بدست می‌آید.

برای مثال، اگر یک سطر از ماتریس C در الگوریتم COUNTMIN را در نظر بگیریم، برابر است با Ax که در ستون i -ام مؤلفه $h(i)$ برابر با ۱ و بقیه مؤلفه‌ها برابر با ۰ است.

۴ پرسش وزین‌های سریع‌تر

یک درخت دودویی کامل با $1 + \log(n)$ طبقه در نظر بگیرید، که هر رأس آن متناظر با یک زیرمجموعه از اندیس‌ها است. به این شکل که هر برگ متناظر با یک اندیس است و هر رأس غیر برگ متناظر با اجتماع دو فرزندش. مثلاً ریشه متناظر با تمام اندیس‌ها است.



عدد متناظر با یک رأس را مجموع x روی اندیس‌های داخل آن رأس در نظر می‌گیریم. به این ترتیب در طبقه j -ام تعداد 2^j عدد در نظر گرفته‌ایم. برای هر طبقه به صورت موازی و مستقل یک الگوریتم COUNTMIN در نظر می‌گیریم، که پارامتر آن $4k$ و احتمال خطای آن $\eta = \delta / 4k \log n$ است. به این ترتیب طبق تحلیل‌های بخش قبل حافظه‌ی استفاده شده از مرتبه‌ی $\mathcal{O}(k \log((k \log n)/\delta) \log n) = \mathcal{O}(k \log(1/\eta) \log(n))$ است.

برای تغییرات، به ازای هر طبقه، تغییر را به رأسی که اندیس مورد نظر را در بر می‌گیرد، اضافه می‌کنیم. همچنین یک شمارنده معمولی نگه می‌داریم که مقدار دقیق $\|x\|_1$ را نگهداری کند.

برای پاسخ دادن پرسش وزین‌ها، دقت کنید که اگر یک اندیس در بردار اصلی وزین باشد (یعنی مقدار بزرگتر از $\|x\|_1/k$ داشته باشد)، تمام رأس‌های روی مسیر از ریشه تا برگ متناظر با این اندیس وزین خواهند بود. به این ترتیب از طبقه‌ی اول شروع می‌کنیم و برای هر طبقه‌ی j یک مجموعه L_j نگه می‌داریم که وزین‌های آن طبقه را در بر داشته باشد. برای شروع L_0 را برابر با ریشه قرار می‌دهیم. برای L_j ، به ازای هر رأس داخل L_{j-1} هر فرزند آن را بررسی می‌کنیم اگر مقدار گزارش شده از $\frac{3}{4}k\|x\|_1$ بیشتر بود، فرزند را در L_j قرار می‌دهیم. مشابه با استدلال‌هایی که در قضیه ۲.۲ آمده بود، از آنجا که COUNTMIN هر طبقه، پارامتر $4k$ دارد، هر نقطه‌ی وزین مقدار گزارش‌شده‌ی بزرگتر از $\frac{3}{4}k\|x\|_1$ دارد. پس طبق خاصیت مسیر ریشه تا برگ که گفته شد، در L_j خواهد آمد. همچنین هر مؤلفه‌ای که مقدار گزارش‌شده‌اش بیش از $\frac{3}{4}k\|x\|_1$ باشد، مقدار واقعی بیش از $\frac{1}{4}k\|x\|_1$ دارد. پس تعداد رأس‌های L_j در هر طبقه حداکثر $2k$ است. توجه کنید اگر پرسش‌های نقطه‌ای همگی به درستی پاسخ داده شوند (یعنی خطایشان قابل قبول داشته باشند)، مجموعه‌ی L_j به ازای هر طبقه، طبق استدلال‌های بالا، تمام وزین‌های طبقه را در بر می‌گیرد و مجموعه نهایی هم تمام اندیس‌های وزین را خواهد داشت. همچنین با فرض درستی پرسش‌های نقطه‌ای در هر طبقه کمتر از $4k$ پرسش می‌کنیم و در کل $Q = 4k \log n$ پرسش می‌کنیم (اگر تعداد پرسش‌هایمان بیش از این شد به این معنا است که یک پاسخ نقطه‌ای اشتباه بوده و می‌توانیم اعلام کنیم که پاسخ را نتوانستیم پیدا کنیم). به این ترتیب این داده ساختار طبق کران اجتماع به احتمال کمتر از $\delta = Q\eta$ ، پاسخ درست می‌دهد که این اثبات درستی الگوریتم را تکمیل می‌کند.

همانطور که گفته شد کمتر از $4k \log n$ پرسش انجام می‌دهیم، پس زمان لازم برای پاسخ دادن پرسش از مرتبه‌ی $O(k \log(1/\eta) \log(n))$ برابر با $O(k \log((k \log n)/\delta) \log n)$ است. همچنین زمان لازم برای تغییر $O(\log n \log((k \log n)/\delta)) = O(\log n \log(1/\eta))$

۵ الگوریتم COUNTSKETCH: پرسش نقطه‌ای برای شمارنده‌ی آزاد و نرم ℓ_2

۱.۵ تعاریف مربوط به ℓ_2

ابتدا در مورد نرم ℓ_2 ، تغییراتی در تعریف ایجاد می‌کنیم. الگوریتم با پارامتر k ، برای پرسش نقطه‌ای باید خطای کمتر از $\|x\|_2/\sqrt{k}$ داشته باشد. همچنین برای پرسش وزین‌ها، مجموعه‌ی خروجی باید $O(k)$ عضو داشته باشد و باید تمام مؤلفه‌هایی که قدرمطلقشان از $\|x\|_2/\sqrt{k}$ بزرگتر است داخل مجموعه قرار باشند.

۲.۵ شرح الگوریتم

الگوریتم COUNTSKETCH شبیه به الگوریتم COUNTMIN است. با کمی تفاوت؛ اول این که به ازای هر سطر به جز تابع درهم‌سازی که هر اندیس از x را به یک اندیس از سطر می‌برد، یک تابع در هم‌ساز 2 - طرفه σ در نظر می‌گیریم که به هر اندیس x عدد ± 1 را اختصاص می‌دهیم. حالا به ازای هر تغییر $x_i \leftarrow x_i + \Delta$ قرار می‌دهیم، $C_{aha(i)} \leftarrow C_{aha(i)} + \sigma_a(i)\Delta$. تفاوت دوم هم این است که برای پرسش نقطه‌ای i ، $\text{median}_a \sigma_a(i) C_{aha(i)}$ را گزارش می‌دهیم. همچنین تعداد ستون‌ها $9k$ قرار داده می‌شود.

۳.۵ اثبات درستی

ابتدا $C_{aha(i)}$ به ازای یک سطر تحلیل می‌کنیم. برای تحلیل پرسش نقطه‌ای i ، طبق تقارن بدون از دست رفتن کلیت مسئله، فرض می‌کنیم $\sigma(i) = 1$. مشابه تحلیل الگوریتم COUNTMIN متغیر تصادفی Z_j را متغیر مشخصه این پدیده تعریف می‌کنیم که $h(i) = h(j)$. به این ترتیب داریم:

$$C_{aha(i)} = x_i + \underbrace{\sum_{j \neq i} \sigma_a(j) Z_j x_j}_E$$

به صورت شهودی مؤلفه‌هایی که σ مثبت دارند، به احتمال مناسبی، مؤلفه‌هایی که σ منفی دارند را «خنثی» می‌کنند. برای اثبات دقیق در نظر داشته باشید که:

$$\mathbb{E}[E] = \sum_{j \neq i} \mathbb{E}[\sigma(j) Z_j x_j] = \sum \mathbb{E}[Z_j x_j \mid \sigma_j] \underbrace{\mathbb{E}[\sigma_j]}_{=0} = 0$$

حالا کافی است یک کران بالا برای $\mathbb{E}[E^2]$ ارائه دهیم، سپس با نامساوی مارکوف نتیجه می‌شود که E به احتمال ثابتی کوچک است:

$$\begin{aligned}\mathbb{E}[E^2] &= \mathbb{E}\left[\left(\sum_{j \neq i} \sigma(j) Z_j x_j\right)^2\right] \\ &= \mathbb{E}\left[\sum_{j \neq i} Z_j x_j^2 + \sum_{j, j' \neq i} \sigma(j) Z_j x_j \sigma(j') Z_{j'} x_{j'}^2\right] \\ &= \sum_{j \neq i} \mathbb{E}[Z_j] x_j^2 + \sum_{j, j' \neq i} \underbrace{\mathbb{E}[\sigma(j)] \mathbb{E}[\sigma(j')]}_{=0} \mathbb{E}[Z_j x_j Z_{j'} x_{j'}^2] \quad (\text{طبق ۲ - طرفه بودن } \sigma) \\ &= \frac{1}{9} \frac{\|x\|_2^2}{k}\end{aligned}$$

حالا طبق نامساوی ینسن داریم:

$$\mathbb{E}[|E|] \leq \sqrt{\mathbb{E}[E^2]} = \frac{1}{3} \frac{\|x\|_2}{\sqrt{k}}$$

و طبق نامساوی مارکوف داریم:

$$\Pr[|E| > \frac{\|x\|_2}{\sqrt{k}}] \leq \frac{1}{3}$$

حالا با اجرای موازی $L = \Theta(\log(1/\delta))$ بار از این الگوریتم، (که در واقع سطرهای ماتریس هستند) و میانگین‌گیری پاسخ‌ها نتیجه می‌شود که

$$\Pr[|\text{median}_a \sigma_a(i) C_{aha(i)} - x_i| > \frac{\|x\|_2}{\sqrt{k}}] < \delta$$

که نتیجه می‌دهد الگوریتم درست کار می‌کند.

۴.۵ تحلیل زمان اجرا و حافظه

مشابه الگوریتم COUNTMIN زمان اجرای لازم برای اعمال تغییر و پرسش نقطه‌ای از مرتبه‌ی $\mathcal{O}(\log(n/\delta))$ است. همچنین زمان لازم برای پرسش وزین‌ها $\mathcal{O}(k \log(n/\delta))$ است. حافظه‌ی لازم هم برای دو نوع پرسش به ترتیب $\mathcal{O}(k \log(n/\delta))$ و $\mathcal{O}(k \log(1/\delta))$ است.

۵.۵ منابع و ارجاع‌ها

- جزوه‌ی درس فشرده‌سازی آقای نلسون [?].
مقاله‌ی مربوط به COUNTMIN و حالت سریع‌تر آن [?].
مقاله‌ی مربوط به COUNTSKETCH [?].