



الگوریتم‌های خلاصه‌سازی برای مه‌داده

محمد هادی فروغ‌مندان‌عربی

پاییز ۱۳۹۹

شمارنده با حافظه کم

جلسه ی سوم

نگارنده: محمد جواد سجادی

۱ مروری بر مباحث گذشته

مسئله ی شمارنده : شمارنده ای داریم که سه عمل انجام می‌دهد :

• `init()`: مقدار عدد n را ۰ می‌گذارد.

• `update()`: عدد n را یک واحد افزایش می‌دهد، یعنی $n \leftarrow n + 1$.

• `query()`: عدد n و یا تخمینی از آن را خروجی می‌دهد.

که میزان حافظه ی دقیق مورد نیاز برای آن $(\log n)$ بررسی شد در این جلسه به راه حل تقریبی با حافظه ی کمتر برای این مسئله می‌پردازیم

۲ راه حل تقریبی

برای تعریف دقیق اینکه منظورمان از تقریب چیست باید گفت که ϵ را اندازه ی خطا و δ را احتمال خطا در نظر می‌گیریم . می‌خواهیم تمامی ϵ و δ را پیدا کنیم که در رابطه ی رو به رو صدق می‌کنند.

$$\mathbb{P}(|\tilde{n} - n| > \epsilon n) < \delta$$

[Mor78]

۱.۲ الگوریتم موریس

اگر بجای اضافه شدن X در هر مرحله اینگونه الگوریتم را در نظر بگیریم :

• $\text{init}()$: مقدار عدد X را ۰ می‌گذارد.

• $\text{update}()$: عدد X را با احتمال 2^{-X} افزایش می‌دهیم.

• $\text{query}()$: عدد $1 - 2^X$ را به عنوان خروجی ارائه می‌دهیم.

اگر X_n را مقدار X پس از n بار به روز رسانی در نظر بگیریم

$$\text{لم } 1) \mathbb{E}[2^{X_n} = n + 1].$$

پایه: حکم به ازای $n=1$ بدهی است. چون پس از ۱ بار به روز رسانی با احتمال $2^0 = 1$ مقدار $X_1 = 1$ می‌شود

$$\mathbb{E}[2^{X_1} = 2^1 = 2 = 1 + 1]$$

گام اثبات : فرض کنید حکم برای $n=k$ درست باشد حکم را برای $n+1$ اثبات می‌کنیم.

$$\begin{aligned} \mathbb{E}[2^{X_{n+1}}] &= \sum_{j=0}^{\infty} \mathbb{E}[2^{X_{n+1}} | X_n = j] * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} (2^{-j} 2^{j+1} + (1 - 2^{-j}) 2^j) * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} (2 + 2^j - 1) * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} (2^j + 1) * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} 2^j * \mathbb{P}(X_n = j) + \sum_{j=0}^{\infty} \mathbb{P}(X_n = j) \end{aligned}$$

و باید گفت که $\sum_{j=0}^{\infty} \mathbb{P}(X_n = j) = 1$ چون تمامی حالات ممکنه برای X_n را در نظر گرفتیم و از طرفی

$$\begin{aligned} \sum_{j=0}^{\infty} 2^j * \mathbb{P}(X_n = j) &= \sum_{j=0}^{\infty} 2^{X_n} * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} 2^{X_n} * \mathbb{P}(X_n = j) = \mathbb{E}[2^{X_n}] \end{aligned}$$

طبق فرض استقرا $\mathbb{E}[2^{X_n}] = n + 1$ پس :

$$\mathbb{E}[2^{X_{n+1}}] + 1 = n + 1 + 1 = n + 2$$

حکم اثبات شد.

$$\text{لم } 2) \mathbb{E}[2^{2 * X_n} = \frac{2}{3} n^2 + \frac{2}{3} n + 1].$$

پایه: حکم به ازای $n=1$ بدهی است. چون پس از ۱ بار به روز رسانی با احتمال $2^0 = 1$ مقدار $X_1 = 1$ می‌شود

$$\mathbb{E}[2^{2 * X_1} = 2^{2 * 1} = 4 = \frac{2}{3} + \frac{2}{3} + 1 = 4]$$

گام اثبات: فرض کنید حکم برای $n=k$ درست باشد حکم را برای $n+1$ اثبات می‌کنیم.

$$\begin{aligned}\mathbb{E}[2^{X_{n+1}}] &= \sum_{j=0}^{\infty} \mathbb{E}[2^{X_{n+1}} | X_n = j] * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} (2^{-j} 2^{j+1} + (1 - 2^{-j}) 2^j) * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} (2^{j+1} + 2^j - 2^j) * \mathbb{P}(X_n = j) \\ &= \sum_{j=0}^{\infty} (2 * 2^j + 2^j) * \mathbb{P}(X_n = j) \\ &= 2 * \sum_{j=0}^{\infty} 2^j * \mathbb{P}(X_n = j) + \sum_{j=0}^{\infty} 2^j * \mathbb{P}(X_n = j)\end{aligned}$$

و باید گفت که $\sum_{j=0}^{\infty} 2^j * \mathbb{P}(X_n = j)$ همان $\mathbb{E}[2^{X_n}]$ است و عبارت $\sum_{j=0}^{\infty} (2^{2*j}) * \mathbb{P}(X_n = j)$ همان $\mathbb{E}[2^{2*X_n}]$ است. پس:

$$\begin{aligned}\mathbb{E}[2^{X_{n+1}}] + 2 * \mathbb{E}[2^{X_n}] &= \frac{2}{\epsilon} n^2 + \frac{2}{\epsilon} n + 1 + 2 * (n + 1) \\ &= \frac{2}{\epsilon} (n^2 + 2n + 1) + \frac{2}{\epsilon} (n + 1) + 1 = \frac{2}{\epsilon} (n + 1)^2 + \frac{2}{\epsilon} (n + 1) + 1\end{aligned}$$

حکم اثبات شد.

می‌خواهیم این مقدار را محاسبه کنیم و یا کرانی برای آن تعیین کنیم

$$\mathbb{P}(|\tilde{n} - n| > \epsilon n)$$

طبق قضیه ی چبیشف خواهیم داشت:

$$\begin{aligned}\mathbb{P}(|\tilde{n} - n| > \epsilon n) &< \frac{1}{\epsilon^2 n^2} \text{Var}[\tilde{n}] \\ &= \frac{\mathbb{E}[\tilde{n}^2] - \mathbb{E}[\tilde{n}]^2}{\epsilon^2 n^2} = \frac{\frac{2}{\epsilon} n^2 + \frac{2}{\epsilon} n + 1 - (n + 1)^2}{\epsilon^2 n^2} \\ &= \frac{\frac{n^2 - n}{\epsilon}}{\epsilon^2 n^2} = \frac{n - 1}{\epsilon n \epsilon^2} < \frac{1}{\epsilon^2}\end{aligned}$$

اما اینجا اشکال بزرگی وجود دارد که مثلاً برای $\epsilon = \frac{1}{2}$ به کران $\frac{1}{\epsilon^2} = 4$ دست پیدا کردیم که اصلاً کران خوبی نیست.

۲.۲ الگوریتم موریس +

به جای قرار دادن یک ماشین با الگوریتم موریس اگر بتوانیم s ماشین با این الگوریتم قرار بدهیم و برای جواب نهایی میانگین خروجی ماشین ها را قرار بدهیم. پس خواهیم داشت:

$$\tilde{n} = \frac{1}{s} \sum_{i=1}^s \tilde{n}_i$$

لم ۳ (ماشین های الگوریتم موریس مستقل هستند). اینگونه میتوان استدلال کرد که برای هر n بار به روز رسانی جواب هر یک مستقل از دیگری خواهد بود به این معنا که اضافه شدن هیچ یک مرتبط به دیگری نیست و تاثیری روی هم ندارند پس این s ماشین مستقل از هم هستند چون همگی تعداد ورودی مشخص و یکسانی ورودی می‌گیرند.

به طریق مشابه نامساوی را می‌نویسیم و با استفاده از لم خواهیم داشت:

$$\begin{aligned} \mathbb{P}(|\tilde{n} - n| > \epsilon n) &< \frac{1}{\epsilon^2 n^2} \text{Var}[\tilde{n}] \\ &= \frac{1}{\epsilon^2 n^2} \text{Var}\left[\frac{1}{s} \sum_{i=1}^s \tilde{n}_i\right] = \frac{1}{\epsilon^2 n^2 s^2} \text{Var}\left[\sum_{i=1}^s \tilde{n}_i\right] \\ &= \frac{1}{\epsilon^2 n^2 s^2} \sum_{i=1}^s \text{Var}[\tilde{n}_i] = \frac{1}{\epsilon^2 n^2 s^2} s * \text{Var}[\tilde{n}_i] = \frac{1}{\epsilon^2 n^2 s} \\ &= \frac{1}{\epsilon^2 n * s} * \frac{n-1}{2} < \frac{1}{2 * \epsilon^2 * s} < \delta \end{aligned}$$

کران کمی بهتر شد اما همچنان به اندازه ی کافی خوب نیست لازم به ذکر است که:

$$s > \frac{1}{2 * \epsilon^2 \delta} = \Theta\left(\frac{1}{\epsilon^2 \delta}\right)$$

۳.۲ الگوریتم موریس++

اجرای t تا موریس+ در کنار هم و با $\delta = \frac{1}{3}$ و به عنوان جواب میانه را بر می‌گردانیم

متغیرهای Y_i را ۱ در نظر می‌گیریم در صورتی که $\frac{1}{4} < \mathbb{P}(|\tilde{n} - n| > \epsilon n)$ باشد در غیر این صورت ۰ در نظر می‌گیریم. پس در صورتی الگوریتم ما بخواهد نتیجه ی مطلوب بدهد پس باید بیش از نیمی از آنها ۱ باشند.

لم ۴ (میانه‌۴). به صورت معمول میانگین را به عنوان خروجی بر می‌گردانند اما در اینجا به علت احتمالاتی بودن سیستم احتمال وجود داده ی پرت ی بسیار زیاد است پس باید تخمین‌گری انتخاب کنیم که نسبت به داده های پرت واکنش زیادی نشان ندهد و یکی از بهترین و ساده ترین این تخمین گر ها میانه است.

لم ۵ $(\mathbb{E}(\sum_{i=1}^t Y_i) > \frac{2t}{3})$. به علت استقلال:

$$\mathbb{E}\left(\sum_{i=1}^t Y_i\right) = \sum_{i=1}^t \mathbb{E}(Y_i) > \sum_{i=1}^t \frac{2}{3} = \frac{2t}{3}$$

حال به بررسی اثبات می‌پردازیم اگر داشته باشیم:

$$\begin{aligned} \mathbb{P}\left(\sum_{i=1}^t Y_i < \frac{t}{2}\right) &= \mathbb{P}\left(\sum_{i=1}^t Y_i - \mathbb{E}[Y_i] < \frac{t}{2} - \mathbb{E}[Y_i]\right) \\ &= \mathbb{P}\left(\mathbb{E}[Y_i] - \sum_{i=1}^t Y_i > \mathbb{E}[Y_i] - \frac{t}{2}\right) = \mathbb{P}\left(|\mathbb{E}[Y_i] - \sum_{i=1}^t Y_i| > \mathbb{E}[Y_i] - \frac{t}{2}\right) \end{aligned}$$

می‌دانیم طبق لم که عبارت سمت راست احتمال مثبت است پس تغییری روی آن ایجاد نمی‌کنیم:

$$\begin{aligned} &\mathbb{P}\left(|\mathbb{E}[Y_i] - \sum_{i=1}^t Y_i| > \mathbb{E}[Y_i] - \frac{t}{2}\right) \\ &= \mathbb{P}\left(|\mathbb{E}[Y_i] - \sum_{i=1}^t Y_i| > \frac{1}{4} \mathbb{E}[Y_i] + \left(\frac{2}{4} \mathbb{E}[Y_i] - \frac{t}{2}\right)\right) \\ &< \mathbb{P}\left(|\mathbb{E}[Y_i] - \sum_{i=1}^t Y_i| > \frac{1}{4} \mathbb{E}[Y_i]\right) \end{aligned}$$

لازم به ذکر است که طبق لم داریم که

$$\frac{2}{4} \mathbb{E}[Y_i] - \frac{t}{2} > \frac{2}{4} * \frac{2t}{3} = \frac{t}{3} - \frac{t}{2} = \frac{t}{6}$$

پس طبق کران چرنوف خواهیم داشت که:

$$\mathbb{P}\left(|\mathbb{E}[Y_i] - \sum_{i=1}^t Y_i| > \frac{1}{4} \mathbb{E}[Y_i]\right) < 2e^{\frac{-e^2 * \mathbb{E}[Y_i]}{2}} < 2e^{\frac{-e^2 * \frac{2t}{3}}{2}} = 2e^{\frac{-\frac{1}{3} * \frac{2t}{3}}{2}} = 2e^{\frac{-t}{9}} < \delta$$

پس با هم هی این تفاسیر خواهیم داشت که :

$$t > c * \ln\left(\frac{1}{\delta}\right) = \Theta\left(\ln \frac{1}{\delta}\right)$$

پس حافظه ی نهایی برابر خواهد بود با :

$$s * t * \ln(\ln(n)) = \Theta\left(\frac{1}{\epsilon^2} * \ln\left(\frac{1}{\delta}\right) * \ln(\ln(n))\right)$$

به عبارت دقیق تر ما به احتمال $1 - \delta$ حافظه ی مورد نیاز ما $O\left(\frac{1}{\epsilon^2} * \ln\left(\frac{1}{\delta}\right) * \ln(\ln\left(\frac{n}{\epsilon\delta}\right))\right)$ خواهد بود.

۴.۲ کمی بهتر؟؟

اگر بجای 2^{-X} عبارت $\frac{1}{1+a}^X$, $a > 0$ مناسبی در جای ان قرار دهیم به کران $O\left(\ln\left(\frac{1}{\epsilon}\right) + \ln\left(*\ln\left(\frac{1}{\delta}\right)\right) + \ln(\ln(n))\right)$ برای حافظه خواهیم رسید و

$$O(\ln(\ln_{1+\epsilon}(n))) = O\left(\ln\left(\frac{1}{\epsilon}\right) + \ln(\ln(n))\right)$$

کران پایین مسئله ی گفته شده است. [PJN17]

مراجع

- [Mor78] Robert Morris. Counting large numbers of events in small registers. *Commun. ACM*, 21(10):840–842, 10 1978.
- [PJN17] Vinh-Kha Le Prof. Jelani Nelson. Sketching algorithms for big data. Lecture 01:3–8, Fall 2017.