



# الگوریتم‌های خلاصه‌سازی برای مه‌داده

محمد هادی فروغمنداعرابی  
پاییز ۱۳۹۹

## خلاصه‌سازی گراف

جلسه هشتم

نگارنده: فاطمه کرمانی

### ۱ خلاصه‌سازی حافظه در الگوریتم پیدا کردن جنگل فراگیر

هدف در این بخش خلاصه‌سازی حافظه‌ی مورد نیاز برای اجرای الگوریتم‌های مختلف روی گراف لازم است. در این جلسه، الگوریتم یافتن جنگل فراگیر روی یک گراف بررسی می‌شود. جنگل فراگیر بزرگ‌ترین زیرگراف یک گراف است که دوری ندارد. مولفه‌های هم‌بندی جنگل فراگیر همان مولفه‌های هم‌بندی گراف است. این مسئله، تعدادی زیرمسئله دارد که در ادامه به آن‌ها می‌پردازیم.

#### ۱.۱ بازیابی بردار $k$ -تنک

تعریف ۱.

به بردار  $x \in \mathbb{R}^n$ ،  $k$ -تنک می‌گوییم اگر  $|\text{support}(x)| < k$ .  $\text{support}(x) \subseteq [n]$  برابر است با مجموعه‌ی اندیس‌های  $i$  که  $x_i \neq 0$ . زیرمسئله‌ی بازیابی  $k$ -تنک برای  $x$ های  $k$ -تنک، مقدار دقیق  $x$  را برمی‌گرداند (اندیس‌های  $\text{support}(x)$  و مقادیرشان) و در بقیه‌ی حالات می‌تواند جواب دل‌خواه دهد.

می‌توان نشان داد این مسئله به صورت قطعی با  $O(k \log(nM))$  بیت حل می‌شود، اگر ورودی‌های مسئله از  $M$  بزرگ‌تر نشوند و مقادیر آپدیت‌ها اعداد صحیح هستند. یک روش معمول برای این کار این است که از یک مپ<sup>۱</sup> استفاده کنیم. در این صورت به اندازه‌ی خانه‌های غیر صفر بیت حافظه نیاز داریم. در این حالت فرض کنید تعداد بیت‌های ناصفر در طول زمان کم و زیاد بشوند. در روش معمول اگر یک بار این تعداد، از  $k$  بیش‌تر شود الگوریتم دیگر درست جواب نخواهد داد حتا اگر بعد از آن با به‌روزرسانی‌های بعدی این تعداد کم‌تر از  $k$  بشود.

<sup>۱</sup>ترجمه‌ش چی می‌شه؟

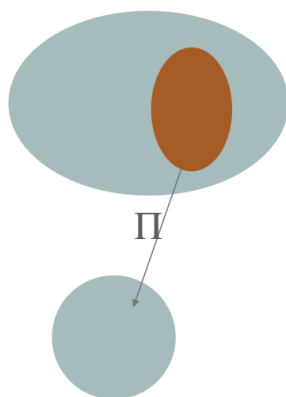


شکل ۱

هدف این است که ساختمان داده ای بسازیم که در مواقعی که در طول انجام به روزرسانی ها هرگاه آرایه ای لحظه ای  $k$ -تنک بود، جواب درست را با احتمال ۱ درست برگرداند و در صورتی که  $k$ -تنک نبود، جواب دل خواه برگرداند. ایده ی طراحی این ساختمان داده این است که به جای خود آرایه ی  $x$ ، یک خلاصه ساز خطی از آن را نگه داریم. این خلاصه ساز را  $\Pi$  می نامیم. به جای  $x$ ،  $\Pi x$  را ذخیره کنیم. ابعاد  $\Pi$ ،  $m \times n$  است و  $m \ll n$ . ویژگی ای که  $\Pi$  باید داشته باشد این است که قابل بازیابی باشد. برای این منظور، اگر در شکل ۴ بیضی بزرگ تر  $\mathbb{R}^n$ ، بیضی کوچک تر  $\mathbb{R}^m$  و بیضی نارنجی رنگ بردارهای  $k$ -تنک باشد،  $\Pi$  هیچ دو عضو نارنجی را به یک موجود مپ نکند. به عبارت دیگر یکی از دو شرط معادل زیر را داشته باشیم:

• برای هر دو  $x, y$   $k$ -تنک،  $\Pi x \neq \Pi y$

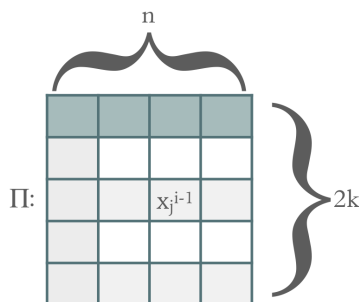
• برای هر  $z$   $2k$ -تنک،  $\Pi z \neq 0$



شکل ۲

توضیح: دو شرط بالا معادلند چون اگر  $x - y$  را برای  $x$  و  $y$  که هر دو  $k$ -تنک هستند را در نظر بگیرد، حاصل یک بردار  $2k$ -تنک است. همچنین شرط دوم به این معنی است که هیچ ترکیب خطی شامل  $2k$  تا از ستون های  $\Pi$  نباید صفر بشود. برای ساختن چنین  $\Pi$  ای یک کاندید مناسب، مطابق شکل ۳ است. در ادامه نشان می دهیم که ماتریس متشکل از انتخاب هر  $2k$  ستون از  $\Pi$ ، ناتبه گون است. نکته: اگر به جای  $\mathbb{R}^n$ ،  $F_p$  که میدان اعداد اول  $p$  تایی است می گرفتیم برای یافتن وارون کار بسیار ساده تری در پیش داریم. در این حالت می توانیم تمام حالت های  $F_p^m$  را بگردیم و ببینیم کدام یکی جواب مورد نظر ماست. پس فرض می کنیم اعدادمان اعضای  $F_p$  هستند و  $p > n$  است و همچنین اعدادی که به عنوان ورودی داریم از  $p$  کوچک تر است.

$x_1, \dots, x_n \in F_p$  در نظر می گیریم و تعریف می کنیم برای هر  $i, j \in [n]$ ،  $i \in [2k]$ ،  $\Pi_{i,j} = x_j^{i-1}$  (یک گزینه ی منطقی برای  $1, 2, \dots, n, x_1, \dots, x_n$  است).



شکل ۳

گزاره ۲. دترمینان ماتریس واندرموند<sup>۲</sup> فرض کنید  $A \in F^{r \times r}$  را تعریف کنیم  $A_{i,j} = x_j^{i-1}$  برای  $i, j \in [r]$  برای یک میدان  $F$ . آنگاه داریم:

$$\det(A) = \prod_{1 \leq i < j \leq r} (x_i - x_j)$$

طبق گزاره ۲ هر  $2k$  ستون از  $\Pi$  را که انتخاب کنیم دترمینان ناصفر دارد.

برای بازیابی  $k$ -تنک، به جای بردار  $x$ ،  $\Pi x$  را ذخیره می‌کنیم. برای به‌روزرسانی، فرض کنیم درایه‌ی  $i$ ام  $x$  را به اندازه  $\Delta$  می‌خواهیم تغییر دهیم. در این صورت حاصل برابر است با  $\Pi x + \Pi^i \Delta$ .  $\Pi^i$  ستون  $i$ ام  $\Pi$  است.

برای بازیابی  $x$ ، چون تعداد بردارهای  $k$ -تنک خیلی زیاد نیست، می‌توانیم همه‌ی آن‌ها را چک کنیم به این صورت که پیدا کنیم که ضرب کدام یک از آن‌ها در ماتریس  $\Pi$  برابر مقدار ذخیره شده هست که آن را به عنوان جواب برگردانیم. این الگوریتم به صورت یکتا جواب درست را به دست می‌آورد چون تبدیل  $\Pi$  یک به یک است.

### تحلیل حافظه

حافظه‌ی مورد نیاز  $2k \log p$  است.

پس زیر مسئله‌ی اول مورد نیاز را بررسی کردیم. در بخش بعد، زیر مسئله‌ی دومی که با آن برخورد می‌کنیم را بیان می‌کنیم.

## ۲.۱ پیدا کردن support

### تعریف مسئله

فرض کنید برای میدان  $F$  می‌خواهیم دو عملیات زیر را با حافظه‌ی کم و به صورت احتمالی داشته باشیم:

•  $\text{update}()$  به روزرسانی  $(\Delta = \pm 1)$

•  $\text{query}(x)$   $i \in \text{support}(x)$  را برگرداند

ایده‌ی اصلی این است که  $\log n$  ساختمان داده در نظر بگیریم و هر خانه‌ی حافظه را به یکی از این ساختمان داده‌ها متناظر کنیم. به طوری که احتمال این که یک خانه‌ی حافظه را به یک ساختمان داده متناظر کنیم به طور نمایی از احتمال این که به ساختمان داده‌ی بعدی متناظر کنیم بیشتر باشد. به عبارت دقیق‌تر اگر تابع تناظر این دو را درهم‌ساز  $h$  در نظر بگیریم داشته باشیم  $\mathbb{P}(h(i) = j) = \frac{1}{p^j}$ . برای این ساختمان داده‌ها یک عدد ثابتی در نظر می‌گیریم که در صورتی که تعداد اعدادی که به آن‌ها متناظر می‌کنیم کمتر از آن عدد ثابت باشد، آن ساختمان داده خوب کار می‌کند. در غیر این صورت عدد دل‌خواه خروجی می‌دهد.

با توجه به این که در هر ساختمان داده، به صورت حدودی نصف قبلی خانه‌ی حافظه وارد می‌شود، یکی از این ساختمان داده‌ها پیدا می‌شود که در بازه‌ی مناسب (تعداد ناصفر اما کمتر از آن عدد ثابت فرض شده) ورودی دارد.

### الگوریتم

• شروع

- تابع درهم‌ساز نمایی  $h$  را تعریف می‌کنیم.  $pr(h(i) = j) = \frac{1}{p^j}$ .

<sup>2</sup>Vandermonde Matrix

–  $A_j$  را برای هر  $j \in [\log n]$  یک الگوریتم بازیابی  $k$  – تنک که در بخش قبل داشتیم قرار می‌دهیم.  
 $k = C \log(1/\delta) *$

• به‌روزرسانی

–  $\text{update}(i, \Delta) = A_{h(i)}.\text{update}(i, \Delta)$

• جواب

– اولین بزرگ‌ترین  $j$  که  $A_j$  خالی نبود: یکی از اعضای  $A_j$ .

1	$A_1$
2	$A_2$
...	...
$\log n$	$A_{\log n}$

شکل ۴

### ۱.۲.۱ تحلیل الگوریتم

در صورتی که  $x = \circ$  در این صورت تمام ساختمان داده‌ها صفر برمی‌گردانند و پاسخ درست است. اگر  $|\text{support}(x)| \leq k$  باشد، در بدترین حالت تمام اعضای  $x$  در یک ساختمان داده قرار می‌گیرند در این حالت نیز جواب درست را برمی‌گردانند.  
پس حالتی را در نظر می‌گیریم که  $k \leq |\text{support}(x)| := t$ .  $T_j$  := تعداد  $x_i$ هایی که در  $A_j$  قرار گرفته‌اند. امید ریاضی  $T_j$  را داریم  $\mathbb{E}T_j = t/2^j$ .  $j^*$  را طوری انتخاب می‌کنیم که داشته باشیم برای  $c$  مناسب  $1 < c < C$ ،  $c \log(1/\delta) \leq t/2^{j^*} < 2c \log(1/\delta)$ ،  $j^*$  با احتمال خوبی اندیس همان ساختمان داده‌ای است که جواب درست را به ما خواهد داد.  
دو اتفاق خوب و تضمین کننده وجود دارد برای این که  $j^*$  پاسخ درست باشد.

$$\mathcal{E}_1 : \max_{j \geq j^*} T_j \leq k \quad \bullet$$

$$\mathcal{E}_2 : T_{j^*} \geq 1 \quad \bullet$$

پس اگر احتمال این که این دو اتفاق بیفتد بالا باشد، در این صورت الگوریتم مطلوب است. احتمال اتفاق نیفتادن  $\mathcal{E}_1$  و  $\mathcal{E}_2$  را حساب می‌کنیم و طبق کران تجمعی<sup>۳</sup> احتمال اتفاق افتادن هر دو را با کران بالای جمع حاصل تخمین می‌زنیم.

$$\mathbb{P}(\neg \mathcal{E}_1) \leq \sum_{j=j^*}^{\log n} \mathbb{P}(T_j > k) = \sum_{j=j^*}^{\log n} \mathbb{P}(T_j > (k 2^j / t) \cdot \mathbb{E}T_j) < (k 2^j / t)^{C'k} \quad (۱)$$

که تساوی با ضرب و تقسیم کردن در  $\mathbb{E} = t/2^j$  برقرار است و نامساوی آخر طبق کران چرنوف با  $1 - e^{-\lambda} > \lambda$  اتفاق می‌افتد:

گزاره ۳. کران چرنوف

$$\mathbb{P}(X > (1 + \lambda)\mu) < \left( \frac{e^\lambda}{(1 + \lambda)^{1 + \lambda}} \right)^\mu < \lambda^{-\Omega(\lambda\mu)} < \sqrt{1/(1 + \lambda)}^{-\Omega(\lambda\mu)}$$

از آنجایی که داشتیم  $k = C \log(1/\delta)$  و  $c \log(1/\delta) \leq t/2^{j^*} \leq 2c \log(1/\delta)$  در نتیجه  $C/2c < k 2^{j^*} / t \leq C/c$  پس در واقع اندازه‌ی  $k 2^{j^*}$  بزرگ است. برای  $\delta$  کوچک و توان به اندازه‌ی کافی بزرگ داریم

$$(k 2^j / t)^{C'k} < (C/2c)^{-C'k} < \delta/2 \quad (۲)$$

<sup>۳</sup>union bound

که با حاصل نامساوی ۱ در کل نتیجه می دهد

$$\mathbb{P}(\neg \mathcal{E}_1) < \delta/2$$

از استدلال های بالا نتیجه می گیریم احتمال این که  $\mathcal{E}$  اتفاق نیفتد زیاد نیست.

حال احتمال رخداد  $\mathcal{E}$  را بررسی می کنیم. می دانیم به طور متوسط چند تا از ورودی ها در ساختمان داده ای  $i$  قرار می گیرد احتمال قرار گیری یکی از ورودی ها در  $j^*$  برابر  $1/2^{j^*}$  است. بنابراین احتمال این که هیچ ورودی ای نداشته باشد برابر است با

$$\mathbb{P}(\neg \mathcal{E}_2) = (1 - 1/2^{j^*})^t < \exp(-t/2^{j^*}) < \delta^c < \delta/2 \quad (3)$$

که مجددا در این نتیجه گیری از  $2c \log(1/\delta) \leq t/2^{j^*} \leq c \log(1/\delta)$  استفاده کردیم و شرط کوچک بودن  $\delta$  به اندازه ی کافی. از ۲ و ۳ و استفاده از کران اجتماع احتمال بد بودن کم تر از  $\delta$  است.

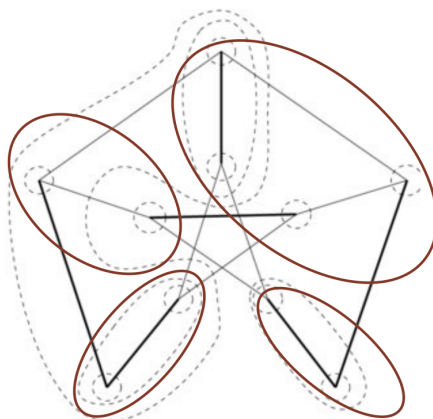
### حافظه ی مورد نیاز

الگوریتم با  $O(\log(1/\delta) \log^2 n)$  بیت حافظه اجرا می شود زیرا  $\log n$  تا بازیابی  $k$ -تنک استفاده می شود که هر کدام یک بردار  $2k$  تایی استفاده می کنند  $2k \simeq \log(1/\delta)$  و حافظه هر کدام  $\log p \simeq \log n$  بیت است.

## ۲ مسئله ی جنگل فراگیر

الگوریتمی برای پیدا کردن یک جنگل فراگیر در یک گراف متغیر <sup>۴</sup> است.

- ورودی: یک گراف
- خروجی: جنگل فراگیر (با احتمال  $(1 - 1/n^b)$ )



شکل ۵

از ساختمان داده ای که در قسمت قبل معرفی کردیم استفاده می کنیم برای حل این مسئله. عملیات ساختمان داده ی مورد نظر:

- اضافه ی یال
- حذف یال
- یک زیرجنگل فراگیر را برگردان

<sup>4</sup>dynamic

## الگوریتم غیر جریانی

قبل از بررسی الگوریتم در حالت کلی، فرض می‌کنیم گراف غیر متغیر است. افزایشی از راس‌ها را ذخیره می‌کنیم. در ابتدا هر راس در کلاس هم‌ارزی جداگانه‌ای قرار داد. هر بار یک یال خروجی از هر مجموعه را انتخاب می‌کنیم و دو مجموعه‌ی دو سر یال‌ها را با هم ادغام می‌کنیم و یال‌های انتخاب شده را در جنگل حاصل قرار می‌دهیم. هر بار تعداد مجموعه‌های افزایش نصف می‌شود. پس بعد از  $\log n$  مرحله یک جنگل فراگیر داریم. برای یک مثال از این الگوریتم شکل ۵

## ۱.۲ روش خلاصه‌سازی AGM

برای هر راس  $u$  یک بردار  $x_u$  تعریف می‌کنیم. این بردار نشان می‌دهد که راس  $u$  به چه راس‌هایی یال دارد.

• اگر یال  $u-v$  وجود داشت آن‌گاه

$$f(x) = \begin{cases} 1, & \text{اگر } u < v \\ -1, & \text{اگر } v < u \end{cases}$$

• اگر یال  $u-v$  وجود نداشت  $x_u[v] = 0$

برای مجموعه‌ی  $A$  که زیر مجموعه‌ای از راس‌های گراف است تعریف می‌کنیم  $x_A := \sum_{u \in A} x_u$ . ویژگی  $x_A$  این است که اگر دو سر یک یال در  $A$  قرار بگیرد، جمع اعداد متناظر دو سر یال با هم ساده می‌شوند. طبق این ویژگی،  $\text{support}(x_A)$  یال‌هایی را مشخص می‌کند که بیرون مجموعه‌ی  $A$  هستند. پس برای هر کدام از مجموعه‌های  $A$ ، با اجرای یک بار الگوریتم  $\text{support-find}(x_A)$  می‌توانیم طبق توضیحات الگوریتم که در بالا ارائه شد، یک یال پیدا کنیم و مجموعه‌های دو سر یال را با هم ادغام کنیم و در صورتی که با اضافه کردن یالی دوری ایجاد نمی‌شد آن یال را به مجموعه‌ی یال‌های جنگل فراگیر مورد نظر اضافه می‌کنیم. و  $x_A$  را به روز رسانی می‌کنیم.

•  $A$  ها را برای شروع تک راس‌های گراف قرار بده.

• برای هر  $x_A$  یک بار  $\text{support-find}(x_A)$  را فراخوانی کن و برای هر مجموعه یک یال پیدا کن. با  $k$  و  $\delta < 1/(Rn^{1+b})$

• برای  $R$

- برای هر مجموعه‌ی  $A$ ، یک یال را پیدا کن.

- برای هر یال: دو مجموعه‌ی  $A$  و  $B$  که دو سر یال هستند را ادغام کن.

$$x_A + x_B = x_{A \cup B}$$

-  $\text{support-find}()$  های دو مجموعه را ادغام کن.

## تحلیل حافظه

برای این الگوریتم  $R$  بار از  $\text{support-find}$  استفاده می‌کنیم. برای ذخیره‌ی اطلاعات مربوط به گراف در  $\text{support-find}$ ،  $n$  بیت حافظه نیاز دارد. در کل این الگوریتم از  $O(nR \log^3 n)$  بیت حافظه استفاده می‌کند.