



جلسه‌ی ۱۴: میانه‌ها و آماره‌های ترتیبی

نگارنده: حسام رهنما-نوید مشایخی

مدرّس: دکتر شهرام خزائی

۱ مقدمه

منظور از i امین آماره‌ی ترتیبی مجموعه‌ای از n عدد، عددی از آن مجموعه است که دقیقاً از $1 - i$ عدد دیگر در آن مجموعه بزرگتر است یا به عبارت دیگر در لیست مرتب شده‌ی (صعودی) این اعداد، در مکان i ام قرار داشته باشد (i امین عضو کوچک مجموعه‌ی آن اعداد باشد). بنابراین به طور مثال اولین آماره‌ی مجموعه‌ای از n عدد همان مینیمم مجموعه‌ی اعداد می‌باشد.

در این بخش از درس، هدف ما پیدا کردن i امین آماره می‌باشد و سعی می‌کنیم الگوریتمی بهینه برای آن ارائه دهیم. پس می‌توانیم مسئله‌ی مطرح شده را به صورت زیر ذکر کنیم:

- ورودی: مجموعه‌ای از n عدد (متمايز) و عدد i که $1 \leq i \leq n$ می‌باشد.
 - خروجی: عددی از مجموعه‌ی اعداد ورودی که i امین عنصر کوچک می‌باشد.
- قبل از حل مسئله‌ی اصلی که پیدا کردن i امین آماره است، به مسئله‌ی دیگری می‌پردازیم.

۲ پیدا کردن همزمان مینیمم و ماکزیمم

می‌خواهیم مقادیر مینیمم و ماکزیمم مجموعه‌ای از n عدد را به طور همزمان پیدا کنیم. ابتدا واضح است که اگر خواهیم هر یک از مقادیر مینیمم و ماکزیمم را به طور جداگانه پیدا کنیم، به راحتی با $n - 1$ مقایسه می‌توان این مقادیر را یافت، پس در نهایت با $2n - 2$ مقایسه به راحتی می‌توانیم هر دو مقدار مینیمم و ماکزیمم را پیدا کنیم. حال می‌خواهیم دو مقدار مینیمم و ماکزیمم مجموعه‌ای از اعداد را با کمتر از $2n - 2$ مقایسه به طور همزمان پیدا کنیم.

در اینجا راه حلی ارائه می‌دهیم که با حداکثر $\lfloor \frac{3n}{3} \rfloor$ مقایسه این دو مقدار را پیدا کنیم.

روش کلی کار به این صورت است که ابتدا اعداد آرایه را به دسته‌های دوتایی کنار هم تقسیم می‌کنیم. سپس به ازای هر یک از این دسته‌های دوتایی، با ۳ مقایسه مقادیر مینیمم و ماکسیمم را که تا به حال به دست آورده‌ایم، آپدیت می‌کنیم. ابتدا دو مقدار که در یک دسته قرار دارند را با هم مقایسه کرده، بعد از آن مقدار کوچکتر این دسته را با مینیمم فعلی (که تا الان به دست آورده‌ایم) و مقدار بزرگتر دسته را با ماکسیمم فعلی مقایسه کرده و مقدار آن‌ها را آپدیت می‌کنیم و در نهایت واضح است که مقادیر مینیمم و ماکسیمم را یافته‌ایم.

اگر بخواهیم تعداد مقایسه‌های انجام شده را حساب کنیم، ما به ازای هر دسته‌ی دوتایی، ۳ مقایسه انجام داده‌ایم. بنابراین در کل $\lfloor \frac{3n}{4} \rfloor$ مقایسه برای پیدا کردن دو مقدار مینیمم و ماکسیمم انجام شده است که کمتر از $2n - 2$ مقایسه‌ای است که در اولین راه حل به آن اشاره شد.

۳ پیدا کردن آماره‌ی k ام

در این بخش راه حلی برای پیدا کردن k امین عضو کوچکتر ارائه می‌دهیم که زمان اجرای آن در بدترین حالت $\Theta(n^2)$ می‌باشد ولی در حالت میانگین در $\Theta(n)$ کار می‌کند. ابتدا این الگوریتم را معرفی می‌کنیم و سعی می‌کنیم آن را بهبود دهیم که در بدترین حالت هم، در زمان خطی اجرا شود.

ایده‌ی کلی استفاده از تابع RANDOMIZED-PARTITION است که در الگوریتم RANDOMIZED-QUICKSORT از آن استفاده کردیم. پس در اینجا الگوریتمی که ارائه می‌دهیم تصادفی می‌باشد. به تابع RANDOMIZED-SELECT که در زیر آمده است توجه کنید.

Algorithm 1 RANDOMIZED-SELECT

```

function RANDOMIZED-SELECT(Array  $A[1..n]$ , Index  $p$ , Index  $r$ , Index  $i$ )
  if  $p = r$  then
    return  $A[p]$ 
   $q \leftarrow$  RANDOMIZED-PARTITION( $A, p, r$ )
   $k \leftarrow q - p + 1$ 
  if  $i = k$  then
    return  $A[q]$ 
  if  $i < k$  then
    return RANDOMIZED-SELECT( $A, p, q - 1, i$ )
  if  $i > k$  then
    return RANDOMIZED-SELECT( $A, q + 1, r, i - k$ )

```

تابع RANDOMIZED-SELECT آرایه‌ی A را به همراه اندیس‌های i, r, p را به عنوان ورودی دریافت می‌کند و k امین عضو کوچک زیرآرایه‌ی $A[p..q]$ را به عنوان خروجی برمی‌گرداند. در اولین خط و در حلقه‌ی این شرط چک می‌شود که اگر آرایه به طول ۱ باشد، همان عنصر آرایه برگردانده شود. در خط سوم با استفاده از تابع RANDOMIZED-PARTITION آرایه‌ی $A[p..r]$ را به دو زیرآرایه‌ی $A[p..q-1]$ و $A[q+1..r]$ تقسیم کردیم که تمام اعداد زیرآرایه‌ی اول از $A[q]$ کوچکتر هستند و تمام اعداد زیرآرایه‌ی دوم از $A[q]$ بیشتر هستند (در اصل عنصر $A[q]$ را به عنوان محور^۱ قرار دادیم). در خط ۴ عدد k را که طول زیرآرایه‌ی $A[p..q]$ می‌باشد را محاسبه می‌کنیم و در خط ۵ چک می‌شود که آیا $A[q]$ که k امین عنصر کوچک آرایه‌ی A می‌باشد، k امین عنصر کوچک A می‌باشد یا نه ($i = k$ است یا نه). اگر جواب این سوال مثبت باشد، همان مقدار $A[q]$ برگردانده می‌شود و در غیر این صورت تعیین می‌کنیم که k امین عضو کوچکتر در کدام یک از زیرآرایه‌های $A[p..q-1]$ و $A[q+1..r]$ قرار دارد و متناسب با آن به صورت بازگشتی تابع RANDOMIZED-SELECT را فراخوانی می‌کنیم.

^۱pivot

۱.۳ بدترین حالت

فرض کنید بخواهیم عنصر مینیمم را پیدا کنیم یعنی تابع را به صورت $\text{RANDOMIZED-SELECT}(A, 1, n, 1)$ فراخوانی کنیم. در این صورت اگر تابع $\text{RANDOMIZED-PARTITION}$ در خط سوم، بزرگترین عدد را به عنوان محور انتخاب کنید، در این صورت مجبوریم در بین $n - 1$ عدد دیگر (به جز عدد ماکزیمم) به دنبال عدد مینیمم باشیم، یعنی تابع $\text{RANDOMIZED-SELECT}(A, 1, n - 1, 1)$ اجرا می‌شود. اگر هر دفعه بزرگترین عدد به عنوان محور انتخاب شود، چون تابع $\text{RANDOMIZED-PARTITION}$ در $\Theta(n)$ اجرا می‌شود، زمان اجرایی نهایی برای پیدا کردن مینیمم با این شرایط $\Theta(n^2)$ خواهد بود. با این حال چون الگوریتم ارائه شده تصادفی است، به ازای هر ورودی داده شده، بدترین حالت با احتمال بسیار کمی اتفاق می‌افتد و زمان اجرای این الگوریتم به طور متوسط $\Theta(n)$ می‌باشد.

۲.۳ محاسبه‌ی زمان اجرای میانگین

فرض کنید زمان اجرای تابع RANDOMIZED-SELECT بر روی آرایه‌ی $A[p..r]$ با n عنصر که یک متغیر تصادفی است را با $T(n)$ نمایش دهیم. می‌خواهیم کرانی برای $E[T(n)]$ که متوسط زمان اجرای الگوریتم است، به دست آوریم. ابتدا متغیر تصادفی شاخص X_k را تعریف می‌کنیم. اگر بعد از اجرای تابع $\text{RANDOMIZED-PARTITION}$ زیرآرایه‌ی $A[p..q]$ دارای دقیقاً k عنصر باشد، X_k را برابر ۱، و در غیر این صورت آن را صفر در نظر می‌گیریم. چون در تابع $\text{RANDOMIZED-PARTITION}$ احتمال اینکه هر یک از اعضای آرایه به عنوان محور انتخاب شوند، برابر است، می‌توانیم نتیجه بگیریم که به ازای $1 \leq k \leq n$ داریم:

$$E[X_k] = \frac{1}{n}$$

حال با توجه به اینکه n امین عنصر کوچک نسبت به عنصر محور ($A[q]$) در کدام سمت باشد، تابع بر روی یکی از زیر آرایه‌های $A[p \dots q - 1]$ و یا $A[q + 1 \dots r]$ اجرا می‌شود. بنابراین می‌توانیم کرانی به صورت زیر برای $T(n)$ بنویسیم (عبارت $O(n)$ مربوط به اجرای تابع RANDOMIZED-PARTITION می‌باشد).

$$\begin{aligned} T(n) &\leq \sum_{k=1}^n X_k (T(\max(k-1, n-k)) + O(n)) \\ &= \sum_{k=1}^n X_k T(\max(k-1, n-k)) + O(n) \end{aligned}$$

پس خواهیم داشت:

$$\begin{aligned} E[T(n)] &\leq E\left[\sum_{k=1}^n X_k T(\max(k-1, n-k)) + O(n)\right] \\ &= \sum_{k=1}^n E[X_k T(\max(k-1, n-k))] + O(n) \\ &= \sum_{k=1}^n E[X_k] \cdot E[T(\max(k-1, n-k))] + O(n) \\ &= \sum_{k=1}^n \frac{1}{n} E[T(\max(k-1, n-k))] + O(n) \end{aligned}$$

در سیگمای بالا تعدادی از جملات دو بار تکرار شده‌اند (با توجه به اینکه n فرد است یا زوج). پس می‌توانیم عبارت بالا را به صورت زیر بازنویسی کنیم.

$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} E[T(k)] + O(n)$$

با توجه به تعریف ثابت می‌کنیم که $E[T(n)] = O(n)$. پس فرض می‌کنیم که به ازای مقادیر کوچکتر از n رابطه‌ی ذکر شده برقرار باشد یعنی عددی ثابت و مثبت مانند c وجود دارد که $E[T(k)] \leq ck$. همچنین برای عبارت $O(n)$ عدد ثابت مثبتی مانند a در نظر می‌گیریم. پس از عبارت بالا خواهیم داشت:

$$\begin{aligned} E[T(n)] &\leq \frac{2}{n} \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} ck + an \\ &= \frac{2c}{n} \left(\sum_{k=1}^{n-1} k - \sum_{k=1}^{\lfloor \frac{n}{2} \rfloor - 1} k \right) + an \\ &\leq \frac{2c}{n} \left(\frac{(n-1)n}{2} - \frac{(\frac{n}{2}-2)(\frac{n}{2}-1)}{2} \right) + an \\ &= c \left(\frac{3n}{4} + \frac{1}{2} - \frac{1}{n} \right) + an \\ &\leq \frac{3cn}{4} + \frac{c}{2} + an \\ &= cn - \left(\frac{cn}{4} - \frac{c}{2} - an \right) \end{aligned}$$

اگر نشان دهیم که عبارت بالا کمتر یا مساوی cn است، اثبات کامل می‌شود. برای این منظور کافی است که عبارت داخل پرانتز، بزرگتر یا مساوی صفر شود. با اعمال این شرط و شرط اضافی $c > 4a$ ، کرانی مانند $n \geq t$ به دست می‌آید که t عدد ثابتی برحسب a و c است. اگر به ازای $cn < t$ را $T(n)$ در نظر بگیریم، مشکل حل می‌شود و اثبات کامل می‌شود. پس توانستیم اثبات کنیم که $E[T(n)] = O(n)$.

۴ الگوریتم

در این بخش الگوریتم قطعی را ارائه می‌دهیم که زمان اجرای آن برای پیدا کردن آماره‌ی i ام در بدترین حالت $O(n)$ می‌باشد.

۱.۴ نحوه‌ی عملکرد الگوریتم

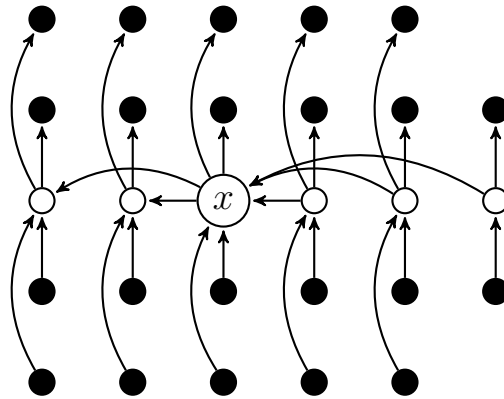
ایده‌ی الگوریتم همانند RANDOMIZED-SELECT می‌باشد با این تفاوت که در سعی می‌کنیم بهترین محور را برای تقسیم آرایه انتخاب کنیم (به جای انتخاب تصادفی). در این الگوریتم از تابع که در آن استفاده کردیم، بهره می‌بریم با این تفاوت که در اینجا تابع عنصر محور را به عنوان ورودی دریافت می‌کند و سپس تقسیم‌بندی را انجام می‌دهد. الگوریتم شامل ۵ مرحله‌ی زیر می‌باشد:

- تقسیم n عدد به دسته‌های ۵ تایی ($\lfloor \frac{n}{5} \rfloor$ دسته‌ی ۵ تایی و حداکثر یک دسته با $n \bmod 5$ عنصر).
- پیدا کردن میانه‌ی هر یک از دسته‌ها (برای این کار می‌توانیم از الگوریتم مرتب‌سازی درجی استفاده کنیم).
- با استفاده از خود تابع میانه‌ی میانه‌هایی که در مرحله‌ی ۲ پیدا کرده‌ایم را به دست می‌آوریم و آن را x می‌نامیم.
- حال از تابع که عنصر محور را به عنوان ورودی دریافت می‌کند، استفاده می‌کنیم. با فراخوانی تابع با ورودی x ، آرایه به دو قسمت تقسیم می‌شود، اعداد سمت چپ x در آرایه‌ی که همگی کوچکتر از x هستند و اعداد سمت راست x که مقادیری بزرگتر از x دارند. فرض کنید x بعد از فراخوانی تابع در جایگاه k ام آرایه باشد (یعنی x عنصر کوچک k ام باشد).
- اگر $i = k$ بود، آنگاه x به عنوان جواب برگردانده می‌شود در غیر این صورت اگر $i < k$ ، باید عنصر کوچک i ام را در زیرآرایه‌ی سمت چپ x پیدا کنیم و اگر $i > k$ باید عنصر کوچک $i - k$ ام را در زیرآرایه‌ی سمت راست پیدا کنیم.

۲.۴ زمان اجرا

برای تحلیل زمان اجرا کران پایینی روی تعداد عناصری که از x بزرگتر هستند، به دست می‌آوریم. شکل زیر دسته‌بندی اعداد و میانه‌ها را مشخص کرده است. در مرحله‌ی دو $\lfloor \frac{n}{5} \rfloor$ میانه به دست آوردیم و در مرحله‌ی ۳، عنصر x را میانه‌ی این $\lfloor \frac{n}{5} \rfloor$ عدد معرفی کردیم، بنابراین حداقل بیش از نصف میانه‌ها، مقداری بزرگتر یا مساوی x دارند و در نتیجه در هر کدام از دسته‌هایی که این میانه‌های بزرگتر از x قرار دارند، حداقل ۳ مقدار بزرگتر از x یافت می‌شود (به جز دسته‌ای که خود x در آن قرار دارد و دسته‌ای که شامل $n \bmod 5$ عدد می‌باشد). پس تعداد اعداد بزرگتر از x حداقل برابر است با:

$$3(\lceil \frac{1}{3} \lceil \frac{n}{5} \rceil \rceil - 2) \geq \frac{3n}{10} - 6$$



این شکل دسته‌های اعداد را مشخص می‌کند که هر دسته‌ی ۵ تایی در یک ستون مشخص شده است و مقدار میانه‌ی هر دسته با دایره‌ی سفید نشان داده شده است. فلش‌ها از اعداد بزرگتر به اعداد کوچکتر رسم شده‌اند و عدد x که میانه‌ی میانه‌ها می‌باشد در شکل مشخص شده است.

به طور مشابه تعداد عناصر کوچکتر از x حداقل برابر $\frac{3n}{10} - 6$ می‌باشد. پس در بدترین حالت، تابع بر روی آرایه‌ای به طول حداکثر $\frac{3n}{10} - 6$ اجرا می‌شود. حال می‌خواهیم رابطه‌ی بازگشتی زمان اجرای $T(n)$ مربوط به تابع را محاسبه کنیم.

در مراحل ۱ و ۲ و ۴ زمان $O(n)$ استفاده شده است (در مرحله‌ی ۲ چون اندازه‌ی هر یک از دسته‌ها ثابت و برابر ۵ است، الگوریتم مرتب‌سازی درجه‌ی برای هر یک از دسته‌ها در $O(1)$ اجرا می‌شود و در کل برای همه‌ی دسته‌ها $O(n)$ خواهد بود). مرحله‌ی ۳ که به طور بازگشتی اجرا می‌شود زمان $T(\lceil \frac{n}{5} \rceil)$ را مصرف می‌کند و در مرحله‌ی ۵ حداکثر زمان $T(\frac{3n}{10} + 6)$ مصرف خواهد شد. بنابراین برای $T(n)$ خواهیم داشت:

$$T(n) \leq T(\frac{n}{5}) + T(\frac{3n}{10} + 6) + O(n)$$

با استفاده از تعریف و فرض وجود ثابت مثبتی مانند c که $T(k) \leq ck$ می‌توان ثابت کرد که رابطه‌ی بازگشتی بالا از $O(n)$ می‌باشد. بنابراین توانستیم الگوریتمی قطعی برای پیدا کردن n امین عضو کوچک مجموعه‌ای از اعداد ارائه دهیم که بدترین زمان اجرای آن $O(n)$ می‌باشد.