

مرتب‌سازی و مرتبه‌ی آماری

n عنصر با کلیدهای a_1, a_2, \dots, a_n و رابطه‌ی ترتیب کامل (total order) \leq بر روی کلیدها داده شده‌اند.

تعریف مرتب‌سازی جای‌گشت π از عناصر داده شده را پیدا کنید به‌طوری‌که:

$$a_{\pi(1)} \leq a_{\pi(2)} \leq \dots \leq a_{\pi(n)}$$

یادآوری

یک رابطه ترتیب جزئی (partial order) است اگر نامتقارن، بازتابی و تراگذری باشد.

رابطه‌ی ترتیب جزئی R بر روی مجموعه‌ی A رابطه‌ی ترتیب کامل است، اگر برای هر $a, b \in A$ داشته باشیم aRb یا bRa .

دسته‌بندی الگوریتم‌های مرتب‌سازی

(۱) از نظر موقعیت داده‌ها در زمان مرتب‌سازی (یعنی داده در کجا ذخیره می‌شود)

◁ داخلی (internal) ▷ خارجی (external)

(۲) از نظر حفظ ترتیب نسبی عناصر پس از مرتب‌سازی

◁ پایدار (stable) ▷ ناپایدار (unstable)

(۳) از نظر نحوه‌ی مرتب‌سازی داده‌ها

◁ مقایسه‌ای (comparison sort) (با مقایسه‌ی عناصر مرتب می‌کند)

◁ غیر مقایسه‌ای (non-comparison sort)

کران پایین الگوریتم‌های مرتب‌سازی

- زمان اجرای هر الگوریتم مرتب‌سازی $\Omega(n)$ است.
- زمان اجرای هر الگوریتم مرتب‌سازی مقایسه‌ای هم در بدترین حالت و هم در حالت میانگین $\Omega(n \lg n)$ است.

یادآوری: تعریف ترتیب الفبایی (Lexicographic Ordering)

فرض:

$$\vec{a} = a_1 a_2 \cdots a_i \cdots a_n$$

$$\vec{b} = b_1 b_2 \cdots b_i \cdots b_m$$

می‌گوییم $\vec{a} \leq \vec{b}$ اگر برای $i < k \leq n$ و $a_i = b_i$ و $a_k < b_k$ یا $n < m$ و برای $1 \leq i \leq n$ داشته باشیم $a_i = b_i$.

مثلاً

$$ab < abc < adc < adda$$

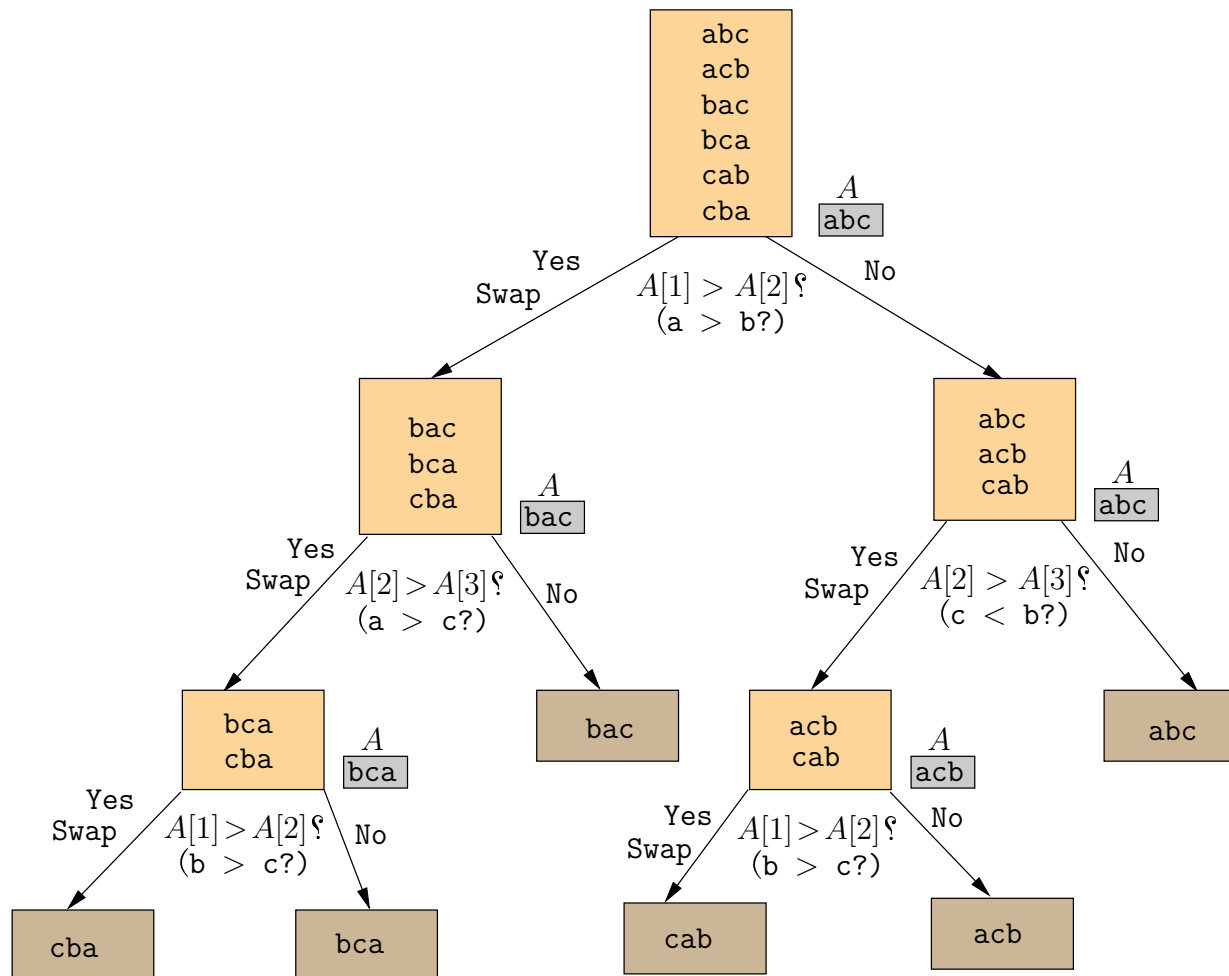
اثبات کران پایین در بدترین حالت

درخت تصمیم (Decision Tree)

هر الگوریتم مقایسه‌ای را می‌توان با یک درخت با ویژگی‌های زیر مدل کرد.

- «حالت مسئله»
 - درخت دودویی کامل است (هر مقایسه یک انشعاب)
 - برگ‌ها حالت نهایی
- مثال: درخت تصمیم برای مرتب‌سازی درجی بر روی $A = [a, b, c]$. هر حالت شامل همه‌ی جای‌گشت‌های ممکن عناصر مرتب است.

داده ساختارها و مبانی الگوریتم‌ها



چندتا لم و قضیه

لم: یک درخت دودویی با ارتفاع h حداکثر 2^h برگ دارد.
این قضیه با استقرا روی h اثبات می‌شود (امتحان کنید).

چند تا لم و قضیه

لم: یک درخت دودویی با ارتفاع h حداکثر 2^h برگ دارد.
این قضیه با استقرا روی h اثبات می‌شود (امتحان کنید).

نکته: حداقل چند تا برگ دارد؟ دودویی کامل و دودویی؟

چندتا لم و قضیه

لم: یک درخت دودویی با ارتفاع h حداکثر 2^h برگ دارد.
این قضیه با استقرا روی h اثبات می‌شود (امتحان کنید).

نکته: حداقل چند تا برگ دارد؟ دودویی کامل و دودویی؟

جواب: ۱ برای درخت عادی و $h + 1$ برای درخت دودویی که هر گره دو یا صفر فرزند داشته باشد.

چندتا لم و قضیه

لم: یک درخت دودویی با ارتفاع h حداکثر 2^h برگ دارد.
این قضیه با استقرا روی h اثبات می‌شود (امتحان کنید).

نتیجه: ارتفاع یک درخت تصمیم که n عنصر را مرتب می‌کند حداقل $\lceil \log n! \rceil$ است.

چندتا لم و قضیه

لم: یک درخت دودویی با ارتفاع h حداکثر 2^h برگ دارد.
این قضیه با استقرا روی h اثبات می‌شود (امتحان کنید).

نتیجه: ارتفاع یک درخت تصمیم که n عنصر را مرتب می‌کند حداقل $\lceil \log n! \rceil$ است.

اثبات: این درخت تصمیم حداقل $n!$ برگ دارد، بنابراین ارتفاعش حداقل $\lceil \log n! \rceil$ است.

چندتا لم و قضیه

لم: یک درخت دودویی با ارتفاع h حداکثر 2^h برگ دارد.
این قضیه با استقرا روی h اثبات می‌شود (امتحان کنید).

نتیجه: ارتفاع یک درخت تصمیم که n عنصر را مرتب می‌کند حداقل $\lceil \log n! \rceil$ است.

اثبات: این درخت تصمیم حداقل $n!$ برگ دارد، بنابراین ارتفاعش حداقل $\lceil \log n! \rceil$ است.

نتیجه: هر الگوریتم مقایسه‌ای که n عنصر را مرتب می‌کند در بدترین حالت حداقل $\lceil \lg n! \rceil$ مقایسه بین عناصر ورودی انجام می‌دهد.

اما می‌دانیم که ..

$n! \leq n^n$ ، پس $\lg n! \leq n \lg n$. ولی این کران بالای ضعیفی است. تقریب استرلینگ (Stirling's Approximation) کران بهتری را می‌دهد. براساس این تقریب داریم:

$$n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \left(1 + \Theta\left(\frac{1}{n}\right)\right)$$

با استفاده از این فرمول می‌توان اثبات کرد که:

$$n! = o(n^n)$$

$$n! = \omega(2^n)$$

$$\lg(n!) = \Theta(n \lg n)$$

کران پایین در حالت میانگین

یعنی چه؟ چه چیز را باید اثبات کنیم؟

کران پایین در حالت میانگین

قضیه: اگر کلیه‌ی جای‌گشت‌های یک ترتیب n تایی با احتمال یک‌سان در ورودی ظاهر شوند، آن‌گاه میانگین عمق برگ‌های درخت تصمیم حداقل $\lg n!$ خواهد بود.

کران پایین در حالت میانگین

اثبات: فرض

- $D(T)$ مجموع عمق برگ‌های یک درخت دودویی T و
 - $D(m)$ کوچک‌ترین مقدار $D(T)$ برای کلیه‌ی درخت‌های دودویی T با m برگ باشد.
- ثابت می‌کنیم $D(m) \geq m \lg m$ ، \Leftarrow میانگین عمق برگ‌های درخت تصمیم $\Omega(\lg n!)$

ادامه‌ی اثبات

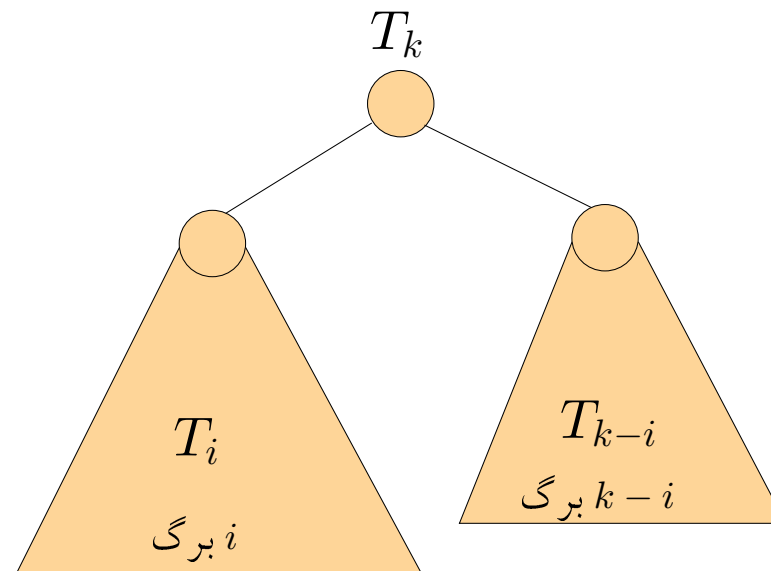
با استقرا ثابت می‌کنیم که $D(m) \geq m \lg m$.

(۱) پایه: $m = 1$ واضح است.

(۲) فرض: برای $m < k$ درست است.

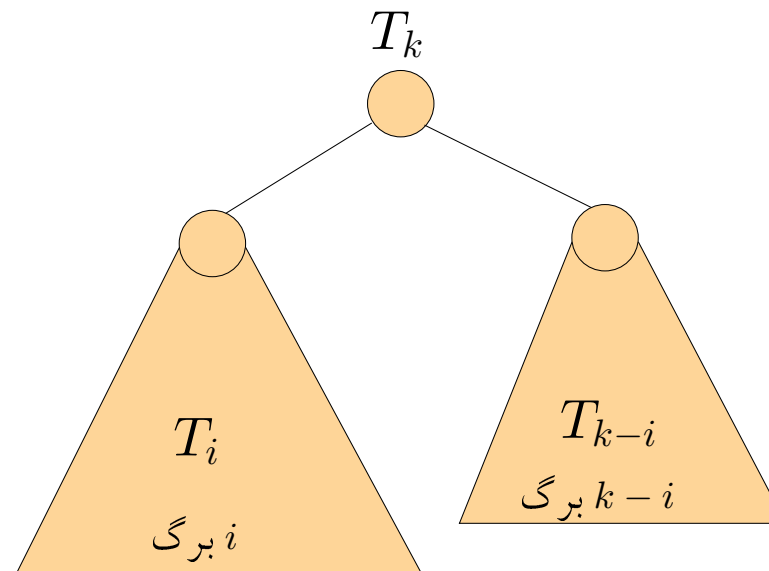
(۳) حکم: T با k برگ را در نظر بگیرید (با T_k نشان می‌دهیم).

داده ساختارها و مبانی الگوریتم‌ها



$$D(T) = i + D(T_i) + (k - i) + D(T_{k-i})$$

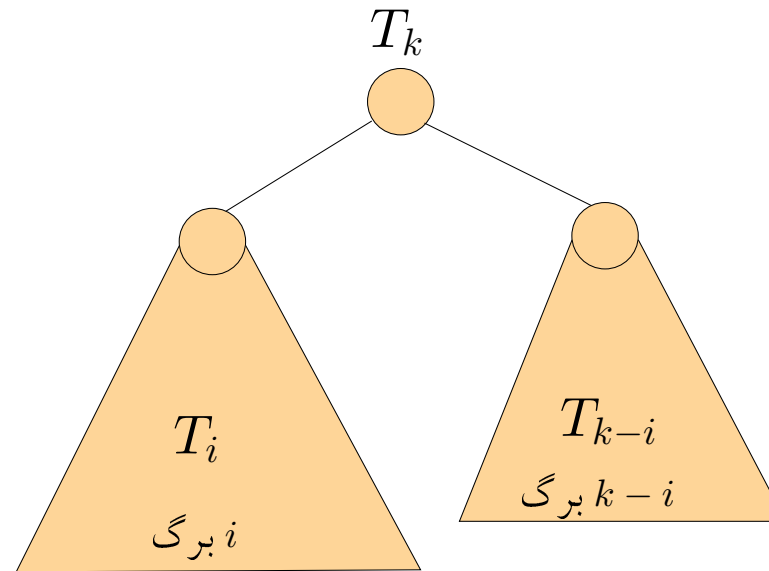
داده ساختارها و مبانی الگوریتم‌ها



$$D(T) = i + D(T_i) + (k - i) + D(T_{k-i})$$

$$D(k) = \min_{1 \leq i \leq k} \{k + D(i) + D(k - i)\}$$

داده ساختارها و مبانی الگوریتم‌ها



$$D(T) = i + D(T_i) + (k - i) + D(T_{k-i})$$

$$D(k) = \min_{1 \leq i \leq k} \{k + D(i) + D(k - i)\}$$

$$D(k) \geq k + \min_{1 \leq i \leq k} \{i \lg i + (k - i) \lg (k - i)\}$$

ادامه‌ی اثبات

برای اعداد طبیعی مقدار کمینه‌ی $i \lg i + (k - i) \lg (k - i)$ در $i = \frac{k}{2}$ اتفاق می‌افتد.

پس

$$D(k) \geq k + k \lg \frac{k}{2} = k \lg k$$

مرتب‌سازی خطی

مرتب‌سازی شمارشی (Count Sort)

ورودی: n عنصر با کلیدهای بین ۱ تا m

COUNT-SORT (A, B, m)

```
1  for  $i \leftarrow 1$  to  $m$ 
2      do  $C[i] \leftarrow 0$ 
3  for  $i \leftarrow 1$  to  $length[A]$ 
4      do  $C[A[i]] \leftarrow C[A[i]] + 1$ 
5  for  $i \leftarrow 2$  to  $m$ 
6      do  $C[i] \leftarrow C[i] + C[i - 1]$ 
7  for  $i \leftarrow length[A]$  downto 1
8      do  $B[C[A[i]]] \leftarrow A[i]$ 
9           $C[A[i]] \leftarrow C[A[i]] - 1$ 
```


چرا درست است؟

چرا درست است؟

- در سطر ۴ تعداد عناصر با کلید یک‌سان را می‌شماریم.
- در پایان حلقه‌ی ۵، $C[i]$ آخرین اندیس آرایه‌ی B است که عناصر با کلید $A[i]$ در آن‌جا قرار می‌گیرند.
- در حلقه‌ی ۷ عناصر در جای خودشان قرار می‌گیرند. این کار به صورت پایدار انجام می‌شود.

تحلیل

- زمان اجرا $O(n)$.
- دلیل آن که این حلقه از انتها به ابتدای A تکرار می‌شود آن است که الگوریتم پایدار شود.
- زمان اجرای کل الگوریتم $O(n + m)$ می‌باشد که اگر m از $O(n)$ باشد زمان اجرای کل $O(n)$ خواهد بود.

حالت خاص

کلیدهای عناصر اعداد ۱ تا n هستند.

COUNT-SORT (A, n)

1 **for** $i \leftarrow 1$ **to** n

2 **do while** $key[A[i]] \neq i$

3 **do** SWAP($A[i], A[key[A[i]]]$)

چرا درست است؟

- الگوریتم به صورت «درجا» مرتب می‌کند.
- هر بار تعویض \leftarrow یک یا دو عنصر در جای نهایی خود.
- از هر کلید بیش از یک عدد \leftarrow الگوریتم ممکن است در حلقه بیفتد.

مرتب‌سازی مبنایی (Radix Sort)

d تعداد رقم‌های اعداد ورودی (رقم i ام بر $i - 1$ ام اولویت دارد)

RADIX-SORT (A, d)

1 **for** $i \leftarrow 1$ **to** d

2 **do** sort array A on digit i by a stable sort

- اگر تعداد رقم‌های ورودی یکسان نباشند؟
- درستی الگوریتم: با استقرا بر روی i
- در انتهای مرحله‌ی i عناصر بر حسب بیت‌ها ۱ تا i کلیدشان مرتب‌اند.
- زمان اجرا به الگوریتم دوم بستگی دارد.
- اگر از مرتب‌سازی شمارشی استفاده کنیم: زمان اجرا $O(dn)$

مرتب‌سازی مبنایی در حالت کلی

- ورودی آرایه‌ای از رکوردها، (یک دسته از کارت‌ها)
- هر رکورد دارای کلیدی با k مؤلفه‌ی $f_1 \dots f_k$ به ترتیب از داده‌گونه‌های $t_1 \dots t_k$
- تعداد مقادیری که هر داده‌گونه t_i می‌تواند داشته باشد محدود و مستقل از n است.

مرتب‌سازی سطلی (Bucket Sort)

BUCKET-SORT (A)

▷ مرتب‌سازی سطلی که لیست A را مرتب می‌کند

```
1  for  $i \leftarrow 1$  to  $k$ 
2      do for each value  $v$  of type  $t_i$ 
3          do make  $B_i[v]$  empty
4      for each record  $r$  in list  $A$  in order
5          do let  $v$  be the value of  $f_i$  of the key for  $r$ 
6              move  $r$  from  $A$  to the end of  $B_i[v]$ 
7      for each value  $v$  of type  $t_i$  from lowest to highest
8          do concatenate  $B_i[v]$  to the end of  $A$ 
```

- ورودی: لیست A با n رکورد، که هر رکورد دارای کلیدی با k مؤلفه به نام‌های $f_1 \dots f_k$ از داده‌گونه‌های $t_1 \dots t_k$
- تعداد حالت‌های t_i برابر s_i
- فرض: f_i بر f_{i-1} اولویت دارد.
- برای $(1 \leq i \leq k)$ داریم: $B_i : \text{array}[t_i]$ of list-type

- اگر هر سطل را به صورت یک صف پیاده‌سازی کنیم، عمل الحاق (concat) در زمان ثابت قابل اجرا است.
- زمان اجرا این الگوریتم برابر است با:

$$\sum_{i=1}^k \Theta(s_i + n) = \Theta(kn + \sum_{i=1}^k s_i)$$

اگر $s_i = \Theta(n)$ آن گاه

$$T(n) = \Theta(kn + \sum_{i=1}^k n) = \Theta(n + kn) = \Theta(n)$$

داده ساختارها و مبانی الگوریتم‌ها

مثال: کلیدها شامل سه مؤلفه با مقادیر $a..z$ ، $۱۰۰...۱$ و $۱۴۰۰...۱۳۰۰$.

عنصر	f_3	f_2	f_1
a_1	a	۵	۱۳۲۰
a_2	c	۱۲	۱۳۱۰
a_3	b	۱۲	۱۳۰۵
a_4	a	۸	۱۴۰۰
a_5	z	۱۰	۱۳۰۸
a_6	b	۱۲	۱۳۰۴
a_7	a	۶	۱۳۱۰

داده ساختارها و مبانی الگوریتم‌ها

عنصر	f_3	f_2	f_1
a_1	a	۵	۱۳۲۰
a_2	c	۱۲	۱۳۱۰
a_3	b	۱۲	۱۳۰۵
a_4	a	۸	۱۴۰۰
a_5	z	۱۰	۱۳۰۸
a_6	b	۱۲	۱۳۰۴
a_7	a	۶	۱۳۱۰

ورودی	$B_1[.]$	خروجی ۱	$B_2[.]$	خروجی ۲	$B_3[.]$	خروجی ۳
a_1	$[۱۳۰۴] a_6$	a_6	$[۵] a_1$	a_1	$['a'] a_1, a_7, a_4$	a_1
a_2	$[۱۳۰۵] a_3$	a_3	$[۶] a_7$	a_7	$['b'] a_6, a_3$	a_7
a_3	$[۱۳۰۸] a_5$	a_5	$[۸] a_4$	a_4	$['c'] a_2$	a_4
a_4	$[۱۳۱۰] a_2, a_7$	a_2	$[۱۰] a_5$	a_5	$['z'] a_5$	a_6
a_5	$[۱۳۲۰] a_1$	a_7	$[۱۲] a_6, a_3, a_2$	a_6		a_3
a_6	$[۱۴۰۰] a_4$	a_1		a_3		a_2
a_7		a_4		a_2		a_5

مرتب‌سازی مقایسه‌ای: مرتب‌سازی سریع

- مرتب‌سازی سریع n عنصر را در بدترین حالت با $O(n)^2$ و در حالت میانگین در $O(n \lg n)$ مرتب می‌کند.
- ضریب ثابت $n \lg n$ کاملاً کوچک است.
- این الگوریتم برای محیط‌های حافظه‌ی خارجی و موازی نیز کارا می‌باشد.

مرتب‌سازی سریع: توصیف الگوریتم

- مرتب‌سازی سریع مانند مرتب‌سازی ادغامی مبتنی بر روش تقسیم و حل است.
- می‌خواهیم آرایه‌ی $A[p..r]$ را مرتب کنیم.
- الگوریتم شامل سه مرحله‌ی زیر است.

(۱) تقسیم: بخش‌بندی (partition) $A[p..r]$ به دو بخش ناتهی $A[p..q]$ و $A[q+1..r]$ ، به طوری که هر عنصر $A[p..q]$ از هر عنصر $A[q+1..r]$ بیش‌تر نباشد.



(۲) حل: دو بخش $A[p..q]$ و $A[q+1..r]$ به صورت بازگشتی مرتب می‌شوند.

(۳) ترکیب: با توجه به ویژگی بخش‌ها، نیازی به ترکیب آن‌ها نیست و کل آرایه مرتب است.

QUICKSORT (A, p, r)

1 **if** $p < r$

2 **then** $q \leftarrow \text{PARTITION}(A, p, r)$

3 QUICKSORT (A, p, q)

4 QUICKSORT($A, q + 1, r$)

به صورت $\text{QUICKSORT}(A, 1, \text{length}[A])$ فراخوانی می شود.

بخش بندی

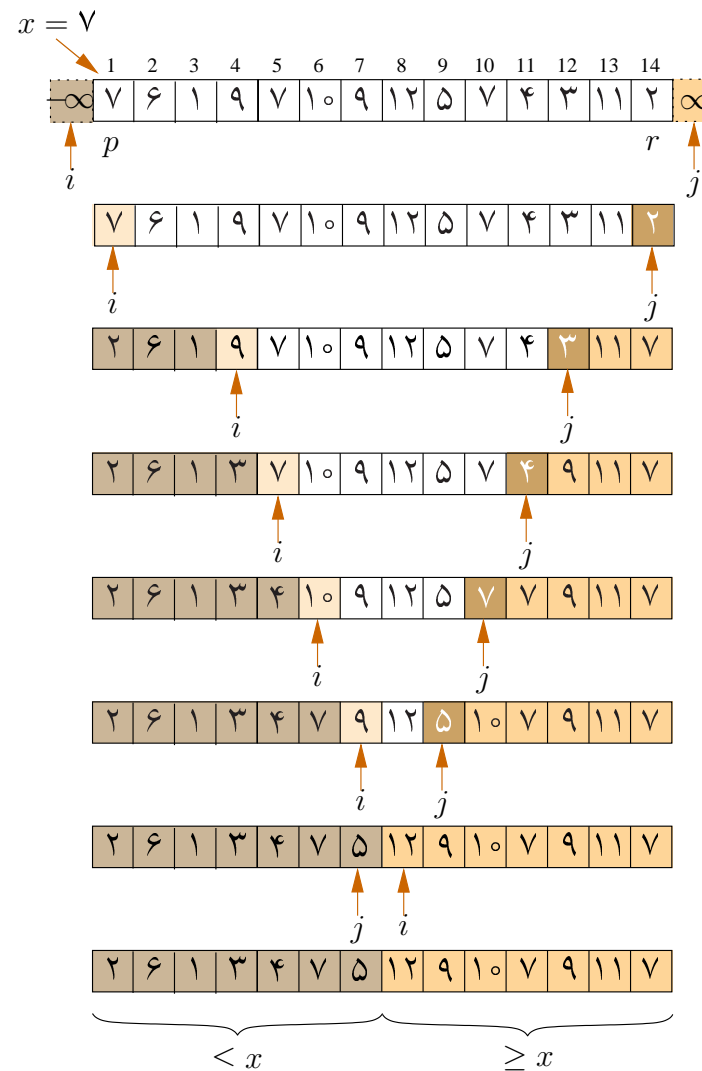
بخش اصلی الگوریتم

PARTITION (A, p, r)

```
1   $x \leftarrow A[p]$ 
2   $i \leftarrow p - 1$ 
3   $j \leftarrow r + 1$ 
4  while true
5      do repeat  $j \leftarrow j - 1$ 
6          until  $A[j] \leq x$ 
7      repeat  $i \leftarrow i + 1$ 
8          until  $A[i] \geq x$ 
9      if  $i < j$ 
10         then  $\text{SWAP}(A[i], A[j])$ 
11         elsereturn  $j$ 
```

مثالی از بخش‌بندی

داده‌ساختارها و مبانی الگوریتم‌ها



چرا بخش‌بندی درست است؟

ویژگی مستقل از حلقه (loop invariant):

$$\bullet \quad i \leq j + 1$$

• به ازای هر مقدار i و j ،

-- بخش $A[p - 1, i - 1]$ حاوی عناصر کم‌تر یا مساوی x

-- بخش $A[j + 1, q + 1]$ حاوی عناصر بیش‌تر یا مساوی محور x

اثبات این ویژگی در انتها، اثبات درستی بخش‌بندی است.

اثبات ویژگی مستقل حلقه

- در ابتدا درست است.
 - فرض: قبل از حلقه‌ی ۴ درست است.
 - j از همه‌ی عناصر بیش‌تر از x عبور می‌کند.
 - i از همه‌ی عناصر کم‌تر از x عبور می‌کند.
 - i نمی‌تواند بیش از یک عنصر از j عبور کند ($i \leq j + 1$)
 - انتهای حلقه، یا $i < j$ ، یا $i = i$ ، یا $i = j + 1$
- که ویژگی برقرار است.

حالت‌های خاص مرتب‌سازی سریع

- آرایه مرتب باشد
 - آرایه برعکس مرتب باشد
 - همه‌ی عناصر آرایه برابر باشد
- تعداد مقایسه‌های کلیدها و تعداد تعویض‌ها را در هر حالت می‌توان دقیقاً شمرد.

مرتب‌سازی سریع (تحلیل الگوریتم)

- بدترین حالت $\Theta(n^2)$ و در حالت میانگین $\Theta(n \log n)$. این کارایی وابسته به نحوه‌ی بخش‌بندی است.
- هزینه‌ی بخش‌بندی برابر $O(n)$ با ثابت کوچک است:
- i و j به عقب بر نمی‌گردند
- تعداد تعویض‌ها حداکثر برابر $n/2$ است.

تحلیل در بدترین حالت

- هنگامی که بخش بندی همیشه n عنصر را به $n - 1$ و 1 عنصر تقسیم کند
- پیچیدگی الگوریتم: $T(n) = T(n - 1) + \Theta(n) = \Theta(n^2)$

تحلیل در به‌ترین حالت

شهود

- بخش‌بندی در هر مرحله n عنصر را به دو آرایه با تعداد عناصر $\lceil \frac{n}{2} \rceil$ و $\lfloor \frac{n}{2} \rfloor$ تقسیم کند.
- پیچیدگی الگوریتم:

$$T(n) = 2T\left(\frac{n}{2}\right) + \Theta(n) \rightarrow T(n) = \Theta(n \log n)$$

بخش‌بندی متوازن

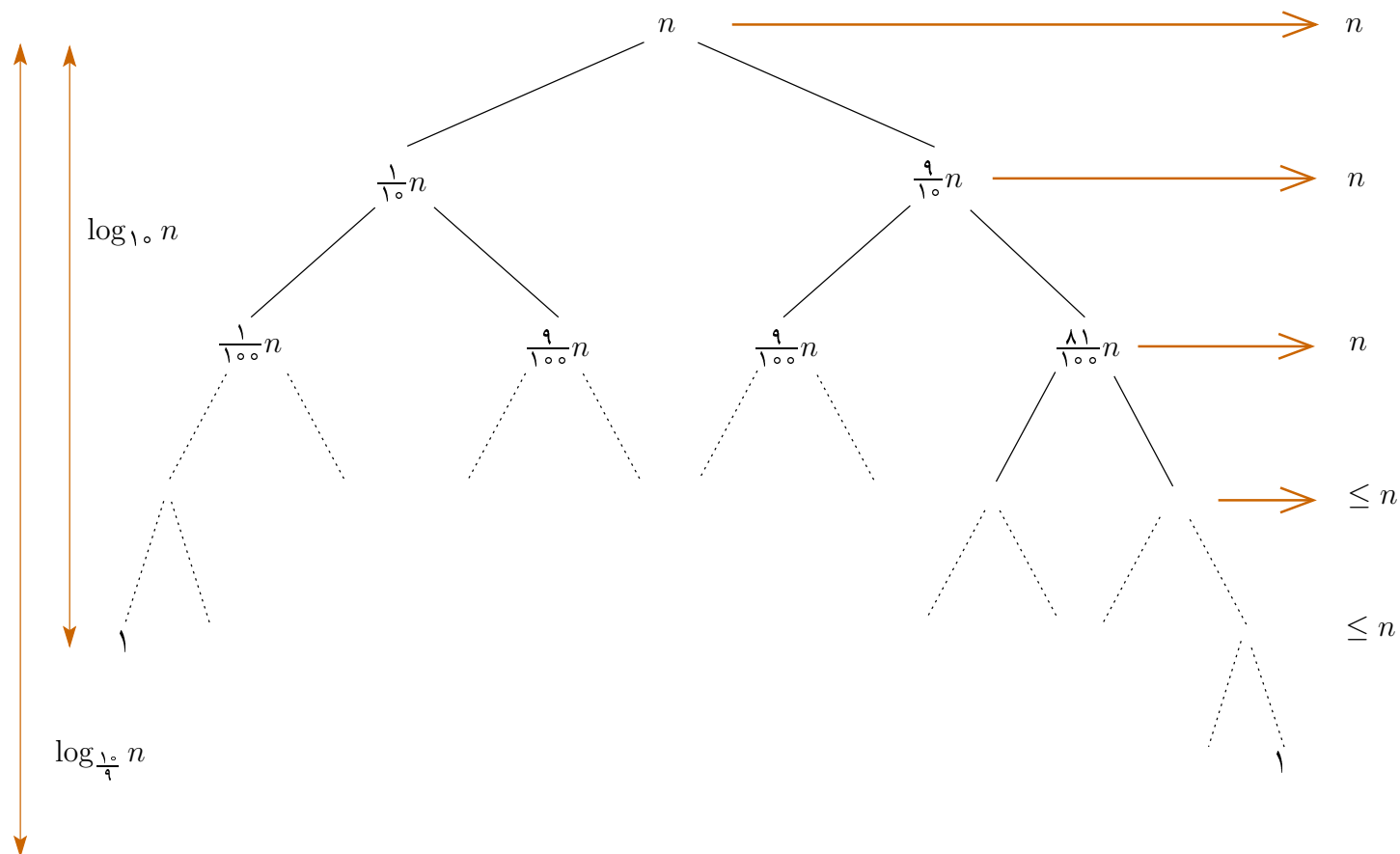
• برای هر ثابت $0 < \alpha < 1$

$$T(n) = T(\alpha n) + T((1 - \alpha)n) + n$$

• مثلاً $T(n) = T(\frac{9}{10}n) + T(\frac{1}{10}n) + n$

• حل $T(n) = \Theta(n \lg n)$ (درخت بازگشت)

درخت بازگشت برای $T(n) = T(\frac{9n}{10}) + T(\frac{n}{10}) + n$



تحلیل در حالت میانگین

برای تحلیل الگوریتم در حالت میانگین، گونه‌ی تصادفی آن را در نظر می‌گیریم.

گونه‌ی تصادفی مرتب‌سازی سریع

به صورت تصادفی یک عنصر دل‌خواه را با عنصر اول تعویض می‌کنیم تا محور شود.

RANDOMIZED-PARTITION (A, p, r)

- 1 $i \leftarrow \text{RANDOM}(p, r)$
- 2 $\text{SWAP}(A[p], A[i])$
- 3 **return** $\text{PARTITION}(A, p, r)$

RANDOMIZED-QUICKSORT (A, p, r)

```
1  if  $p < r$ 
2    then  $q \leftarrow \text{RANDOMIZED-PARTITION}(A, p, r)$ 
3         RANDOMIZED-QUICKSORT( $A, p, q$ )
4         RANDOMIZED-QUICKSORT( $A, q + 1, r$ )
```

تحلیل Randomized-Quicksort در بدترین حالت

- فرض: محور طوری انتخاب شود که آرایه n عضوی به دو بخش q عضوی و $n - q$ عضوی تقسیم می‌شود ($1 \leq q \leq n - 1$).
- در بدترین حالت،

$$T(n) = \max_{1 \leq q \leq n-1} \{T(q) + T(n-q)\} + \Theta(n)$$

- اثبات می‌کنیم که $T(n) = \Theta(n^2)$

ادامه‌ی تحلیل در بدترین حالت

$$T(n) = \max_{1 \leq q \leq n-1} \{T(q) + T(n-q)\} + \Theta(n)$$

ابتدا با استقرای اثبات می‌کنیم که $T(n) = O(n^2)$.

- پایه: برای مقدار ثابت بدیهی است

- فرض: برای $m < n$ ، $T(m) \leq cm^2$

- پس

$$T(n) \leq \max_{1 \leq q \leq n-1} \{cq^2 + c(n-q)^2\} + \Theta(n)$$

داده‌ساختارها و مبانی الگوریتم‌ها

$$T(n) \leq \max_{1 \leq q \leq n-1} \{cq^2 + c(n-q)^2\} + \Theta(n)$$

- می‌دانیم که بیشینه‌ی $q^2 + (n-q)^2$ روی هر دو مرز $q = 1$ و $q = n-1$ اتفاق می‌افتد (با رسم نمودار تابع و با در نظر گرفتن $1 \leq q \leq n-1$)
- پس در حالت $q = 1$ داریم:

$$T(n) \leq c(1 + (n-1)^2) + \Theta(n) = cn^2 - 2c(n-1) + \Theta(n)$$

$$T(n) \leq cn^2$$

به شرطی که c را آنقدر بزرگ بگیریم که $\Theta(n)$ را پوشاند.

ادامه‌ی تحلیل در بدترین حالت

$$T(n) = \max_{1 \leq q \leq n-1} \{T(q) + T(n-q)\} + \Theta(n)$$

نشان می‌دهیم که $T(n) = \Omega(n^2)$.

• فرض: $T(m) \geq cm^2$ برای $m < n$

• پس $T(n) \geq \max_{1 \leq q \leq n-1} \{cq^2 + c(n-q)^2\} + \Theta(n)$

• بنابراین، $T(n) \geq cn^2 - 2c(n-1) + \Theta(n) \geq cn^2$

به شرطی که c آن قدر کوچک باشد که $\Theta(n) - 2c(n-1)$ مثبت شود.

داده ساختارها و مبانی الگوریتم‌ها

• پس $T(n) = \Theta(n^2) \iff T(n) = \Omega(n^2) = O(n^2)$

تحلیل مرتب‌سازی سریع تصادفی در حالت میانگین

- فرض می‌کنیم محور با احتمال یکسان از بین n عنصر انتخاب می‌شود.
- احتمال این که مرتبه‌ی محور k باشد برابر $\frac{1}{n}$ است.
- طبق این الگوریتم بخش‌بندی، در هر دو حالت $k = 1$ و $k = 2$ ، آرایه‌ی A به دو بخش با اندازه‌های 1 و $n - 1$ تقسیم می‌شود.

• تعریف: $\bar{T}(n)$: میانگین زمان اجرا

$$\bar{T}(n) = \frac{1}{n} \left[\bar{T}(1) + \bar{T}(n-1) + \sum_{q=1}^{n-1} [\bar{T}(q) + \bar{T}(n-q)] \right] + \Theta(n)$$

• می‌دانیم که در بدترین حالت $T(1) = \Theta(1)$, $T(n-1) = \Theta(n^2)$ و $\bar{T}(n) \leq T(n)$

• پس:

$$\frac{1}{n} [\bar{T}(1) + \bar{T}(n-1)] \leq \frac{1}{n} [\Theta(1) + \Theta(n^2)] = \Theta(n)$$

• بنابراین

$$\bar{T}(n) = \frac{1}{n} \sum_{q=1}^{n-1} [\bar{T}(q) + \bar{T}(n-q)] + \Theta(n) = \frac{2}{n} \sum_{k=1}^{n-1} \bar{T}(k) + \Theta(n)$$

$$\bar{T}(n) = \frac{1}{n} \sum_{k=1}^{n-1} \bar{T}(k) + \Theta(n) \quad \text{حل}$$

• حدس می‌زنیم $\bar{T}(n) \leq an \lg n + b$ برای $a, b > 0$

$$\begin{aligned} \bar{T}(n) &= \frac{1}{n} \sum_{k=1}^{n-1} \bar{T}(k) + \Theta(n) \\ &\leq \frac{1}{n} \sum_{k=1}^{n-1} [ak \lg k + b] + \Theta(n) \\ &= \frac{1}{n} \sum_{k=1}^{n-1} ak \lg k + \frac{1}{n} \sum_{k=1}^{n-1} b + \Theta(n) \\ &= \frac{1}{n} \sum_{k=1}^{n-1} k \lg k + \frac{1}{n} b(n-1) + \Theta(n) \end{aligned}$$

و از تقریب انتگرال ثابت می‌شود که:

$$\sum_{k=1}^{n-1} k \lg k \leq \int_{x=1}^{n-1} x \lg x = \frac{1}{2} n^2 \lg n - \frac{1}{8} n^2$$

پس

$$\begin{aligned} \bar{T}(n) &\leq \frac{2a}{n} \left(\frac{n^2}{2} \lg n - \frac{n^2}{8} \right) + \frac{2b}{n} (n-1) + \Theta(n) \\ &\leq an \lg n + b + [\Theta(n) + b - \frac{an}{4}] \end{aligned}$$

a و b را طوری انتخاب می‌کنیم تا درون پرانتز منفی شود.

$$\bar{T}(n) \leq an \lg n + b$$

بنابراین

$$\bar{T}(n) = O(n \lg n) \text{ یا}$$

مرتب‌سازی با کم‌ترین تعداد مقایسه برای تعداد عناصر کم

- حداقل تعداد مقایسه‌ی مورد نیاز (بین عناصر) برای مرتب‌سازی n عنصر $\lceil \lg n! \rceil$ است.
- برای n های کوچک آیا الگوریتمی با کم‌ترین تعداد مقایسه‌ی ممکن وجود دارد؟
- برای n برابر ۱، ۲، ۳ مرتب‌سازی «درج دودویی» (binary insertion sort) تعداد مقایسه‌ی بهینه دارد (برابر ۰، ۱، ۳)
- ولی برای $n > ۳$ بهینه نیست.

مرتب‌سازی درج دودویی

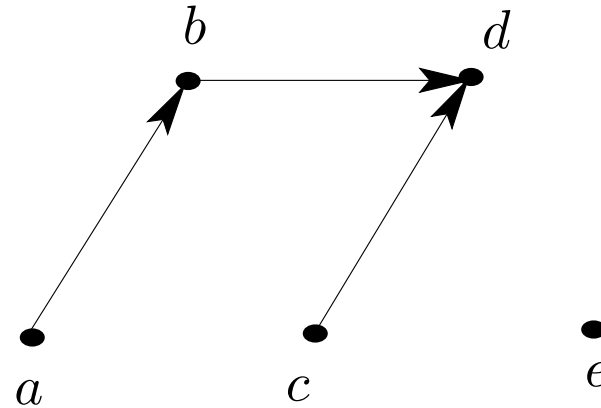
اگر $B(n)$ تعداد دقیق مقایسه‌های الگوریتم درج دودویی (در بدترین حالت ورودی) برای مرتب‌سازی n عنصر باشد، مقادیر $B(n)$ برای n های کوچک و مقایسه‌ی آن با $\lceil \lg n! \rceil$ در جدول زیر آمده است.

$n =$	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰	۱۱	۱۲	۱۳	۱۴	۱۵	۱۶	۱۷
$\lceil \lg n! \rceil =$	۰	۱	۳	۵	۷	۱۰	۱۳	۱۶	۱۹	۲۲	۲۶	۲۹	۳۳	۳۷	۴۱	۴۵	۴۹
$B(n) =$	۰	۱	۳	۵	۸	۱۱	۱۴	۱۷	۲۱	۲۵	۲۹	۳۳	۳۷	۴۱	۴۵	۴۹	۵۴

مرتب‌سازی بهینه‌ی ۵ عنصر

آیا الگوریتمی با ۷ مقایسه برای مرتب‌سازی ۵ عنصر (در بدترین حالت ورودی) وجود دارد؟

داده ساختارها و مبانی الگوریتم‌ها



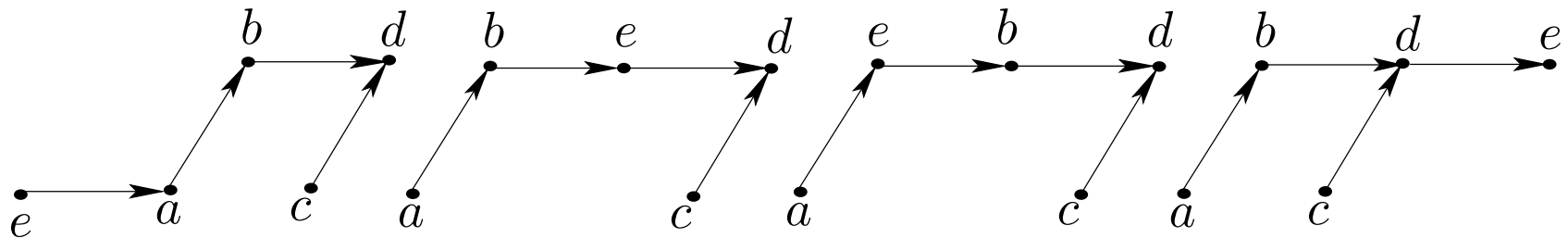
پس از ۳ مقایسه

a با b (فرض b بزرگ‌تر است)

d با c (فرض c بزرگ‌تر است)

b با c (فرض c بزرگ‌تر است)

ابتدا e را در زنجیره‌ی $a \rightarrow b \rightarrow c$ (با ۲ مقایسه در بدترین حالت) درج می‌کنیم.
حالات مختلف پس از درج e



با حداکثر ۳ مقایسه‌ی دیگر d در زنجیره درج می‌شود.

اگر این عناصر K_1 تا K_5 باشند، روش این کار مطابق زیر است:

(۱) K_1 را با K_2 و K_3 را با K_4 مقایسه کن و عناصر کوچک‌تر و بزرگ‌تر را پیدا کن.

(۲) با یک مقایسه‌ی دیگر دو عنصر کوچک بند فوق را با هم مقایسه کن. در انتهای این مرحله ۵ عنصر به صورت شکل در می‌آید. اگر این عناصر a تا e باشند، داریم:
$$a \leq b \leq d \text{ و } c \leq d$$

(۳) e را در زنجیره‌ی $a < b < d$ به روش دودویی درج کن، این کار حداکثر ۲ مقایسه نیاز دارد و یکی از حالات شکل قبل حاصل می‌شود.

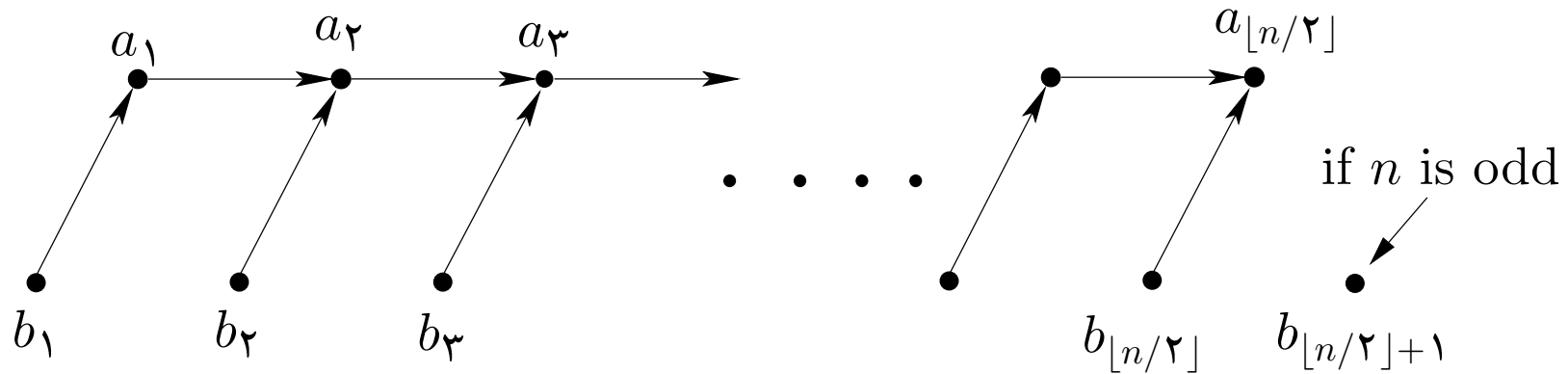
(۴) با توجه به این که $c < d$ ، درج c در زنجیره‌ی چهارتایی از عناصری که در مراحل بالا مرتب شده‌اند، به روش دودویی حداکثر ۲ مقایسه دیگر نیاز دارد.

بنابراین این کار حداکثر ۷ مقایسه نیاز دارد.

تعمیم الگوریتم (الگوریتم فورد-جانسون)

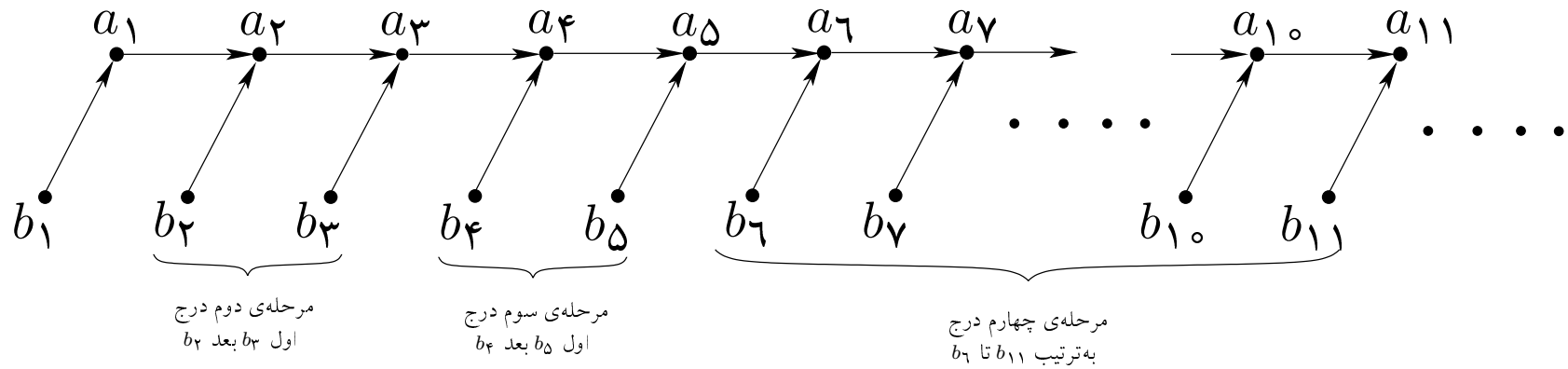
الگوریتم بالا به طرز جالبی توسط آقایان فورد (Lester Ford Jr.) و جانسون (Slemer Johnson) در سال ۱۹۵۹ تعمیم داده شد. ایده‌ی اصلی این الگوریتم تشکیل یک «زنجیره‌ی اصلی» از عناصری است که مرتب هستند و تعیین یک ترتیب مشخص برای درج دیگر عناصر در این زنجیره به طوری که در انتها زنجیره‌ی اصلی شامل همه‌ی عناصر شود.

داده ساختارها و مبانی الگوریتم‌ها



نحوی درج در زنجیره‌ی اصلی در مرحله‌ی k ام

داده ساختارها و مبانی الگوریتم‌ها



نحوه‌ی درج در زنجیره‌ی اصلی در مرحله‌ی k ام.

$$2t_{k-1} + (t_k - t_{k-1}) - 1 = t_k + t_{k-1} - 1 = 2^k - 1$$

الگوریتم برای مرتب‌سازی n عنصر به صورت زیر است:

- (۱) عناصر را به $\lfloor \frac{n}{2} \rfloor$ زوج‌عنصر و احتمالاً یک عنصر اضافی (اگر n فر باشد) تقسیم کن.
 - (۲) هر زوج‌عنصر را با هم مقایسه کن و آن‌ها را مرتب کن.
 - (۳) به صورت بازگشتی $\lfloor \frac{n}{2} \rfloor$ عنصر کوچک‌تر (از زوج‌عناصر) را مرتب کن و پس از آن عناصر را مطابق شکل نام‌گذاری کن. عناصر $a_1 \leq a_2 \leq \dots \leq a_{\lfloor \frac{n}{2} \rfloor}$ را «زنجیره‌ی اصلی» می‌نامیم.
 - (۴) عناصر b در شکل را با ترتیب زیر به بخش اولیه‌ی زنجیره‌ی اصلی به روش دودویی درج کن. این ترتیب طوری انتخاب شده است که تعداد عناصر بخش اول زنجیره‌ی اصلی $2^k - 1$ باشد تا با k مقایسه بتوان زنجیره‌ی کامل را ایجاد کرد:
- (a) ابتدا b_3 را به زنجیره‌ی $a_1 \leq a_2$ و سپس b_2 را به زنجیره‌ی $a_1 \leq b_1$ و احتمالاً b_3 درج کن. تعداد عناصر زنجیره حداکثر ۳ و درج با ۲ مقایسه امکان‌پذیر است.

داده‌ساختارها و مبانی الگوریتم‌ها

(b) ابتدا b_5 را به زنجیره‌ی شامل ۷ عنصر a_1 تا a_4 و b_1 تا b_3 درج کن. سپس b_4 را به زنجیره‌ی a_1 تا a_3 ، b_1 تا b_3 و احتمالاً b_4 درج کن. تعداد عناصر زنجیره حداکثر ۷ و هر درج با ۳ مقایسه امکان‌پذیر است.

(c) این روند را ادامه بده. در حالت کلی برای $t_1, t_2, \dots = 1, 3, 5, 11, \dots$ در مرحله‌ی k ام هر عنصر b_{t_k} تا $b_{t_{k-1}+1}$ را به ترتیب در زنجیره‌ی عناصری که حتماً از آن عنصر کوچک‌ترند درج کن. آخرین مرحله شامل عنصر $b_{\lfloor \frac{n}{4} \rfloor + 1}$ (در صورت وجود) می‌شود.

محاسبات و تحلیل الگوریتم

t_k ها را طوری تعیین کنیم که

$$2t_{k-1} + (t_k - t_{k-1}) - 1 = t_k + t_{k-1} - 1 = 2^k - 1$$

می‌دانیم که $t_0 = 1$

پس کافی است رابطه‌ی بازگشتی $t_k = 2^k - t_{k-1}$ را برای $k > 0$ حل کنیم.
با بازکردن و جای‌گذاری روشن است که

$$t_k = \sum_{i=0}^k (-1)^{k-i} 2^i$$

با استفاده از

$$2t_k = \sum_{i=0}^k (-1)^{k-i} 2^{i+2} = \sum_{i=2}^{k+2} (-1)^{k-i} 2^i$$

داریم

$$\begin{aligned} 2t_k - t_k &= 2^{k+2} - 2^{k+1} - (-1)^{k-1} 2^1 - (-1)^k 2^0 \\ &= 2^{k+1} - [2(-1)^{k-1} + (-1)^k] \\ &= 2^{k+1} - (-1)^{k-1} \\ &= 2^{k+1} + (-1)^k \end{aligned}$$

و از آن داریم:

$$t_k = \frac{1}{3} [2^{k+1} + (-1)^k]$$

تحلیل الگوریتم

$F(n)$ تعداد مقایسه‌ها

$$F(n) = \lfloor \frac{n}{2} \rfloor + F(\lfloor \frac{n}{2} \rfloor) + G(\lceil \frac{n}{2} \rceil)$$

$G(\lceil \frac{n}{2} \rceil) =$ حداکثر تعداد مقایسه‌های لازم برای درج عناصر b در زنجیره‌ی اصلی

اگر $t_{k-1} \leq m < t_k$

می‌دانیم که درج هر عنصر b_k که $t_{j-1} \leq k \leq t_j$ حداکثر j مقایسه نیاز دارد. پس

$$G(m) = \sum_{j=1}^{k-1} j(t_j - t_{j-1}) + k(m - t_{k-1})$$

و داریم

$$\begin{aligned} \sum_{j=1}^{k-1} j(t_j - t_{j-1}) &= [(k-1)t_{k-1} - (k-1)t_{k-2}] + \\ &\quad [(k-2)t_{k-2} - (k-2)t_{k-3}] + \dots + (t_1 - t_0) \\ &= (k-1)t_{k-1} - t_{k-2} - t_{k-3} - \dots - t_1 - t_0. \end{aligned}$$

و بنابراین

$$G(m) = km - (t_{k-1} + t_{k-2} + \dots + t_1 + t_0) \leq km$$

داده ساختارها و مبانی الگوریتم‌ها

برای یک مقدار $k \geq 1$ داریم

$$m \geq t_{k-1} = (2^k + (-1)^{k-1})/3 \geq 2^{k-2}$$

پس $k \leq \lg m + 2$

بنابراین

$$G(m) \leq km \leq m \lg m + 2m$$

پس

$$F(n) \leq \frac{n}{2} + F(n/2) + \lceil \frac{n}{2} \rceil \lg \lceil \frac{n}{2} \rceil + 2 \lceil \frac{n}{2} \rceil \leq F(\frac{n}{2}) + \frac{n}{2} \lg n + O(n)$$

که نتیجه می‌گیریم

$$n \lg n \leq F(n) \leq n \lg n + O(n)$$

برای محاسبه‌ی دقیق‌تر، فرض کنید

$$w_k = t_0 + t_1 + \dots + t_{k-1} = \lfloor \frac{2^{k+1}}{3} \rfloor$$

در این صورت،

$$(w_0, w_1, w_2, w_3, w_4, \dots) = (0, 1, 2, 5, 10, 21, \dots)$$

می‌توان اثبات کرد که

$$F(n) - F(n-1) = k \iff w_k < n \leq w_{k+1}$$

و شرط آخر معادل است با

$$\frac{2^{k+1}}{3} < n < \frac{2^{k+2}}{3} \implies k+1 < \lg(3n) \leq k+2 \implies k = \lg\left(\frac{3}{4}n\right)$$

بنابراین

$$F(n) - F(n-1) = \lceil \lg\left(\frac{3}{4}n\right) \rceil$$

و جواب این رابطه‌ی بازگشتی به صورت زیر است:

$$F(n) = \sum_{k=1}^n \left\lceil \lg\left(\frac{3}{4}k\right) \right\rceil$$

کارایی الگوریتم فورد-جانسون

$n =$	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰	۱۱	۱۲	۱۳	۱۴	۱۵	۱۶	۱۷
$\lceil \lg n! \rceil =$	۰	۱	۳	۵	۷	۱۰	۱۳	۱۶	۱۹	۲۲	۲۶	۲۹	۳۳	۳۷	۴۱	۴۵	۴۹
$F(n) =$	۰	۱	۳	۵	۷	۱۰	۱۳	۱۶	۱۹	۲۲	۲۶	۳۰	۳۴	۳۸	۴۲	۴۶	۵۰

$n =$	۱۸	۱۹	۲۰	۲۱	۲۲	۲۳	۲۴	۲۵	۲۶	۲۷	۲۸	۲۹	۳۰	۳۱
$\lceil \lg n! \rceil =$	۵۳	۵۷	۶۲	۶۶	۷۰	۷۵	۸۰	۸۴	۸۹	۹۴	۹۸	۱۰۳	۱۰۸	۱۱۳
$F(n) =$	۵۴	۵۸	۶۲	۶۶	۷۱	۷۶	۸۱	۸۶	۹۱	۹۶	۱۰۱	۱۰۶	۱۱۱	۱۱۶

جدول ۱: تعداد مقایسه‌های الگوریتم فورد-جانسون $F(n)$ برای تعداد کم عناصر و مقایسه‌ی آن با کران پایین.

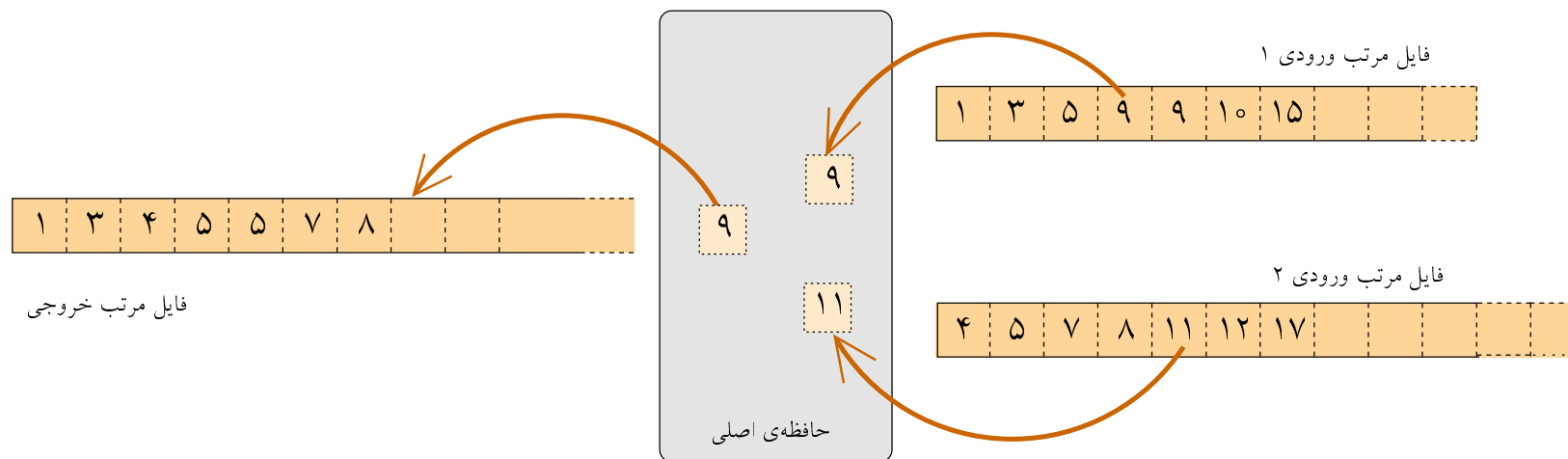
مرتب‌سازی خارجی (فرض)

- اطلاعات بر روی فایل‌ها به صورت ترتیبی ذخیره شده است.
- هر فایل شامل n رکورد است. هر رکورد یک کلید دارد.
- می‌خواهیم در فایل خروجی رکوردها براساس کلیدهایشان مرتب باشند.
- با هر دسترسی به دیسک k رکورد خوانده می‌شود.
- تعداد فایل‌هایی که در یک زمان باز هستند r و محدود است.
- تعداد حافظه‌ی اصلی قابل استفاده ثابت است.
- عملیات مقایسه و محاسبات فقط می‌تواند در حافظه‌ی اصلی انجام شود.

معیار کارایی: تعداد دسترسی‌ها به دیسک

مرتب‌سازی خارجی ادغامی (External Merge Sort)

ادغام دو قطعه‌ی مرتب



داده‌ساختارها و مبانی الگوریتم‌ها

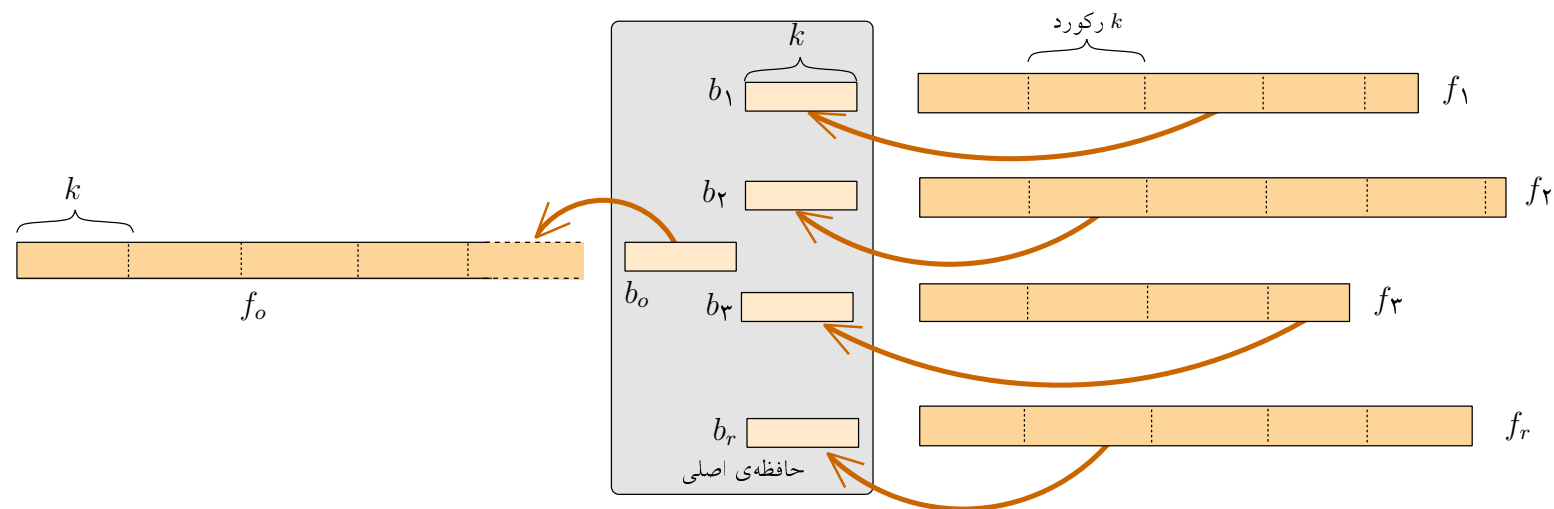
قطعه‌ی اول n_1 رکورد و قطعه‌ی دوم n_2 رکورد

$k = 1$ ، با $n_1 + n_2$ بار خواندن و همین تعداد نوشتن می‌توان ادغام را انجام داد.

ولی در حالت کلی با $\lceil \frac{n_1}{k} \rceil + \lceil \frac{n_2}{k} \rceil$ بار.

مقدار حافظه‌ی مورد نیاز: به‌اندازه‌ی $3k$ رکورد.

ادغام چند فایل مرتب

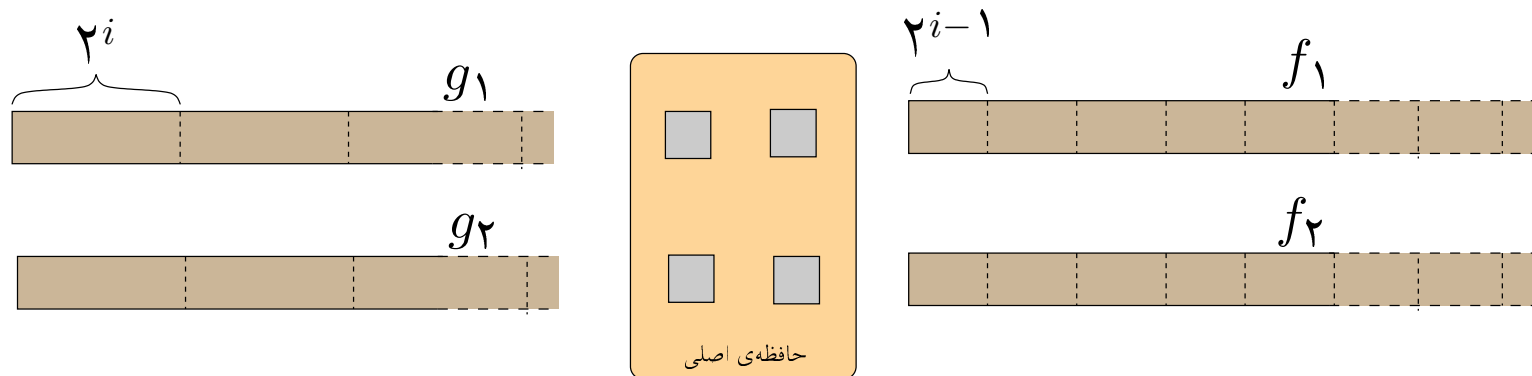


f_i به اندازه‌ی n_i رکورد دارد.

تعداد دسترسی‌ها: $\sum_{i=1}^r \lceil \frac{n_i}{k} \rceil$ بار.

مقدار حافظه‌ی مورد نیاز: به اندازه‌ی $(r + 1)k$ رکورد.

مرتب‌سازی خارجی ادغامی



مرتب‌سازی خارجی ادغامی (الگوریتم کلی)

برای $k = 1$ بیان می‌شود. چهار فایل f_1, f_2, g_1 و g_2 احتیاج است.

(۱) فایل ورودی را به دو فایل f_1 و f_2 با حداکثر تعداد یک رکورد اختلاف تقسیم کن.

(۲) برای $i = 1, \dots, M$ مراحل زیر را تکرار کن:

در این مرحله فرض می‌کنیم که f_1 و f_2 (یا g_1 و g_2) شامل قطعاتی به طول 2^{i-1} هستند و هر قطعه مرتب است و تعداد قطعات دو فایل ورودی حداکثر یک واحد اختلاف دارد.

داده ساختارها و مبانی الگوریتم‌ها

۲-۱) f_1 و f_2 را به صورت فایل‌های ورودی در نظر می‌گیریم. قطعات با شماره‌های یک‌سان f_1 و f_2 را با یکدیگر ادغام کن و قطعه‌ای به طول دو برابر ایجاد کن. حاصل این ادغام قطعاتی مرتب به طول 2^i (به جز حداکثر یک قطعه به طول کم‌تر) است این قطعات را به ترتیب یک‌بار در g_1 و بار دیگر در g_2 بنویس.

۲-۲) g_1 و g_2 را به عنوان فایل‌های ورودی و f_1 و f_2 را به عنوان فایل‌های خروجی در نظر بگیر و مرحله‌ی بالا را تکرار کن.

داده‌ساختارها و مبانی الگوریتم‌ها

مثال مرتب‌سازی خارجی ادغامی

f_{in}	۲۸ ۳ ۹۳ ۱۰ ۵۴ ۶۵ ۳۰ ۹۰ ۱۰ ۶۹ ۸ ۲۲ ۳۱ ۵ ۹۶ ۴۰ ۸۵ ۹ ۳۹ ۱۳ ۸ ۷۷ ۱۰																							
f_1	۲۸ ۳ ۹۳ ۱۰ ۵۴ ۶۵ ۳۰ ۹۰ ۱۰ ۶۹ ۸ ۲۲																							
f_2	۳۱ ۵ ۹۶ ۴۰ ۸۵ ۹ ۳۹ ۱۳ ۸ ۷۷ ۱۰																							
g_1	۲۸ ۳۱ ۹۳ ۹۶ ۵۴ ۸۵ ۳۰ ۳۹ ۸ ۱۰ ۸ ۱۰																							
g_2	۳ ۵ ۱۰ ۴۰ ۹ ۶۵ ۱۳ ۹۰ ۶۹ ۷۷ ۲۲																							
f_1	۳ ۵ ۲۸ ۳۱ ۹ ۵۴ ۶۵ ۸۵ ۸ ۱۰ ۶۹ ۷۷																							
f_2	۱۰ ۴۰ ۹۳ ۹۶ ۱۳ ۳۰ ۳۹ ۹۰ ۸ ۱۰ ۲۲																							
g_1	۳ ۵ ۱۰ ۲۸ ۳۱ ۴۰ ۹۳ ۹۶ ۸ ۸ ۱۰ ۱۰ ۲۲ ۶۹ ۷۷																							
g_2	۹ ۱۳ ۳۰ ۳۹ ۵۴ ۶۵ ۸۵ ۹۰																							
f_1	۳ ۵ ۹ ۱۰ ۱۳ ۲۸ ۳۰ ۳۱ ۳۹ ۴۰ ۵۴ ۶۵ ۸۵ ۹۰ ۹۳ ۹۶																							
f_2	۸ ۸ ۱۰ ۱۰ ۲۲ ۶۹ ۷۷																							
	۳ ۵ ۸ ۸ ۹ ۱۰ ۱۰ ۱۰ ۱۳ ۲۲ ۲۸ ۳۰ ۳۱ ۳۹ ۴۰ ۵۴ ۶۵ ۶۹ ۷۷ ۸۵ ۹۰ ۹۳ ۹۶																							
	$g_1 = f_{out}$																							

درستی و تحلیل

- با استقرا می‌توان نشان داد که در انتهای مرحله‌ی i هر فایل خروجی دارای قطعه‌هایی مرتب و به طول 2^i است، به‌جز حداکثر یک قطعه که طولش از 2^i کم‌تر است. هم‌چنین تعداد قطعه‌های دو فایل خروجی حداکثر یک واحد اختلاف دارند.
- بنابراین برای $M = \lceil \log n \rceil$ (تعداد تکرار حلقه) یکی از فایل‌های خروجی حاوی یک قطعه‌ی مرتب شامل تمام n رکورد فایل ورودی و دیگری خالی است.
- با توجه به این که در هر تکرار همه‌ی n رکورد یک‌بار خوانده و یک‌بار نوشته می‌شوند، تعداد دسترسی به دیسک در مجموع برابر $2n(\lceil \log n \rceil + 1)$ است ($2n$ بار خواندن و نوشتن برای تقسیم فایل اصلی).

داده‌ساختارها و مبانی الگوریتم‌ها

• برای حالت $k > 1$ ، این تعداد برابر $(\lceil \log n \rceil + 1) 2^{\lceil \frac{n}{k} \rceil}$ خواهد بود.

- با تقسیم فایل ورودی به r فایل با اندازه‌هایی یکسان و با استفاده از r حافظه نیز می‌توان فایل را مرتب کرد در آن صورت، تعداد دست‌رسی به دیسک $2n(\lceil \log_r n \rceil + 1)$ می‌شود و در حالت کلی برابر $2\lceil \frac{n}{k} \rceil (\lceil \log_r n \rceil + 1)$.
- در حالت کلی به حافظه‌ای به اندازه‌ی $2rk$ نیاز است.
- حالت کلی را مرتب‌ساز ادغامی چندگانه (Multiway Merge) می‌گوییم.

مرتب‌سازی خارجی چندفازه (Polyphase)

مثال: $n = ۳۴$

پس از گام	$f_۱$	$f_۲$	$f_۳$
ابتدا	۱۳(۱)	۲۱(۱)	—
۱	—	۸(۱)	۱۳(۲)
۲	۸(۳)	—	۵(۲)
۳	۳(۳)	۵(۵)	—
۴	—	۲(۵)	۳(۸)
۵	۲(۱۳)	—	۱(۸)
۶	۱(۱۳)	۱(۲۱)	—
۷	—	—	۱(۳۴)

مرتب‌سازی خارجی چندفازه

- سه (در حالت کلی به $r + 1$) فایل لازم داریم.
- به فایل ورودی کم‌ترین تعداد رکورد با کلید ∞ اضافه کن تا n برابر F_i (امین عدد فیبانوچی) شود.
- فایل ورودی را به دو فایل با اندازه‌های F_{i-1} و F_{i-2} تقسیم کن ($F_i = F_{i-1} + F_{i-2}$)

• باندازه‌ی $M = ?$ بار تکرار کن

-- از سه فایل f_1 دارای F_m قطعه‌ی مرتب (در ابتدا $F_m = F_i$) است که اندازهی هر قطعه‌ی آن F_r است (در ابتدا $F_r = F_o = 1$).

هم‌چنین f_2 دارای F_{m+1} قطعه‌ی مرتب است که اندازهی هر قطعه‌ی آن F_{r+1} است. و f_3 خالی است.

-- F_m قطعه‌ی مرتب از فایل f_1 را با همین تعداد قطعه‌ی مرتب از فایل f_2 را با هم ادغام کن و حاصل را در f_3 بنویس.

-- حال f_1 خالی، f_2 دارای $F_{m-1} = F_{m+1} - F_m$ قطعه به اندازهی F_{r+1} و f_3 دارای F_m قطعه‌ی مرتب به اندازهی $F_{r+2} = F_r + F_{r+1}$ است.

-- نام‌گذاری فایل‌ها را به تناسب تغییر بده.

اثبات درستی و مقدار M

• رابطه‌ی $a(b) + c(d) \longrightarrow e(f) + g(h)$

یعنی یک فایل حاوی a قطعه‌ی مرتب هر کدام به‌اندازه‌ی b ، با فایلی که c قطعه‌ی مرتب هر کدام به‌اندازه‌ی d دارد ادغام می‌شوند و دو فایل حاوی e و f قطعه‌ی مرتب به‌اندازه‌های f و h ایجاد می‌کنند.

• ویژگی مستقل از حلقه:

-- در هر مرحله اعداد a تا h فوق اعداد فیبوناچی هستند.

-- در ابتدا و انتهای هر مرحله یکی از فایل‌ها تهی است.

-- اگر F_i کوچک‌ترین عدد فیبوناچی بزرگ‌تر یا مساوی n باشد، در گام $i - r - 1$ (برای $3 \leq r \leq i - 2$) از الگوریتم، رابطه‌ی زیر برقرار است:

$$F_r(F_{i-r+1}) + F_{r-1}(F_{i-r}) \longrightarrow F_{r-1}(F_{i-r+2}) + F_{r-2}(F_{i-r+1})$$

• پایه:

$$F_{i-1}(F_2) + F_{i-2}(F_1) \longrightarrow F_{i-2}(F_3) + F_{i-3}(F_2)$$

- چون $F_{i-2} < F_{i-1}$ ، به تعداد F_{i-2} رکورد از دو فایل ورودی می‌خوانیم و در یکی از فایل‌های خروجی می‌نویسیم.
- تعداد قطعات این فایل خروجی برابر F_{i-2} و اندازه‌ی هر قطعه‌ی آن برابر $F_1 + F_2 = F_3$ خواهد بود.
- از فایل بزرگ‌تر به تعداد $F_{i-1} - F_{i-2} = F_{i-3}$ قطعه به اندازه‌ی F_2 باقی می‌ماند.
- با فرض درستی پایه‌ی استقرا، می‌توان فرمول بالا را ثابت کرد.

- در انتهای گام $i - 4$ ام (برای $r = 3$) دو فایل خروجی هر کدام یک قطعه دارند و اندازه‌هایشان به ترتیب برابر F_{i-1} و F_{i-2} است.
- پس، در انتهای گام $i - 3$ ام همه چیز در یک فایل مرتب می‌شود.
- تعداد گام‌های الگوریتم برابر $i - 4$ است و در گام $i - r - 1$ ام الگوریتم به تعداد ${}^2F_{i-2}F_{i-r+2}$ رکورد می‌خواند و می‌نویسد.
- با احتساب تقسیم اولیه‌ی فایل که به اندازه‌ی 2F_i دسترسی به دیسک نیاز دارد، در مجموع الگوریتم به

$${}^2F_i + \sum_{r=3}^{i-2} {}^2F_{i-2}F_{i-r+2}$$

حافظه‌ی خارجی دسترسی خواهد داشت.

مرتب‌ی آماری

کوچک‌ترین و بزرگ‌ترین عنصر

با $\lceil \frac{3n}{4} - 2 \rceil$ مقایسه

دو کوچک‌ترین عناصر

راه حل بدیهی: با $n - 1 + n - 2$ مقایسه

راه حل به تر:

- ایجاد یک درخت مقایسه با n برگ و ارتفاع $\lceil \lg n \rceil$ برای تعیین کوچک‌ترین عنصر در ریشه (که دقیقاً $n - 1$ مقایسه نیاز دارد، برابر تعداد گره‌های داخلی این درخت با n برگ)

- دومین کوچک‌ترین عنصر حتماً در این درخت با اولین عنصر مقایسه شده است.
- پس دومین کوچک‌ترین عنصر، کوچک‌ترین عنصر بین $\lceil \lg n \rceil$ عنصر است
- با $\lceil \lg n \rceil - 1$ مقایسه دومین کوچک‌ترین عنصر یافت می‌شود.
- در مجموع $n - 1 + \lceil \lg n \rceil - 1$ مقایسه.

تعمیم راه حل قبل برای ۳ کوچک‌ترین عناصر:

- سومین کوچک‌ترین عنصر یا با دومین و یا اولین کوچک‌ترین عنصر مقایسه شده است.
- اولین کوچک‌ترین عنصر حداکثر با $\lceil \lg n \rceil - 1$ عنصر (غیر از کوچک‌ترین عنصر) مقایسه شده است.
- دومین عنصر حداکثر با $\lceil \lg n \rceil - 2$ عنصر (غیر از کوچک‌ترین و دومین کوچک‌ترین عنصر) مقایسه شده است.
- پس باید کوچک‌ترین عنصر بین حداکثر $\lceil \lg n \rceil - 3$ عنصر را پیدا کنیم که حداکثر $\lceil \lg n \rceil - 4$ مقایسه نیاز داد
- پس در مجموع حداکثر $n - 1 + (\lceil \lg n \rceil - 1) + \lceil \lg n \rceil - 3 = n + \lceil \lg n \rceil - 4$ مقایسه

داده ساختارها و مبانی الگوریتم‌ها

نیاز است..

تعمیم برای $k > 3$

- با روش گذشته با $n - 1 + \sum_{i=2}^k i \lceil \lg n \rceil - (k + 1) = n + \mathcal{O}(k^2 \lg n)$ مقایسه
- با استفاده از هرم کمینه: در $\mathcal{O}(n + k \lg n)$
- با یافتن k امین عنصر (x) و بخش‌بندی حول x در $\mathcal{O}(n + k \lg k)$

یافتن k امین عنصر

(۱) با میانگین $O(n)$ براساس بخش‌بندی در مرتب‌سازی سریع

(۲) $O(n)$ در بدترین حالت

براساس بخش‌بندی ولی با تضمین این که اندازه‌های دو بخش $O(n)$ هستند

$O(n)$ k امین عنصر با میانگین

RANDOMIZED-SELECT ((A, p, r, i))

▷ Find the i th element in $A[p..r]$, assuming $1 \leq i \leq r - p + 1$

1 **if** $p = r$

2 **then return** $A[p]$

3 $q \leftarrow$ RANDOMIZED-PARTITION (A, p, r)

4 $k \leftarrow q - p + 1$

5 **if** $i \leq k$

6 **then return** RANDOMIZED-SELECT (A, p, q, i)

7 **elsereturn** RANDOMIZED-SELECT ($A, q + 1, r, i - k$)

تحلیل Randomized-Select

- فرض می‌کنیم که عناصر نامساوی هستند (حالت بدتر)
- اگر تعداد عناصر n باشد، محور با احتمال $\frac{1}{n}$ ، مرتبه‌اش $1 \leq k \leq n$ است.
- اگر $k = 1, 2$ باشد، آرایه به دو بخش 1 عضوی و $n - 1$ عضوی تقسیم می‌شود.
- در حالت کلی n ، آرایه به دو بخش k عضوی و $n - k$ عضوی تقسیم می‌شود.

پس

$$T(n) \leq \frac{1}{n} \left[T(\max(1, n-1)) + \sum_{k=1}^{n-1} T(\max(k, n-k)) \right] + O(n)$$

$$T(n) \leq \frac{1}{n} \left[T(\max(1, n-1)) + \sum_{k=1}^{n-1} T(\max(k, n-k)) \right] + O(n)$$

داریم

$$\max\{k, n-k\} = \begin{cases} k & \text{if } k \geq \lceil n/2 \rceil, \\ n-k & \text{if } k < \lceil n/2 \rceil \end{cases}$$

- n فرد: جمله‌های $T(\lceil n/2 \rceil + 1), \dots, T(n-1)$ دوبار در Σ ظاهر می‌شوند،
 $T(\lceil n/2 \rceil)$
- n زوج: جمله‌های $T(\lceil n/2 \rceil + 1), T(\lceil n/2 \rceil + 2), \dots, T(n-1)$ دوبار و جمله‌ی $T(\lceil n/2 \rceil)$ یک بار ظاهر می‌شوند.
- جمله‌ی $\frac{1}{n}T(n-1)$ را می‌توان در مقابل $O(n)$ نادیده گرفت.

$$\begin{aligned}
 T(n) &\leq \frac{1}{n} \left[T(\max(1, n-1)) + \sum_{k=1}^{n-1} T(\max(k, n-k)) \right] + O(n) \\
 &\leq \frac{1}{n} \left[T(n-1) + 2 \sum_{k=\lceil n/2 \rceil}^{n-1} T(k) \right] + O(n) \\
 &= \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} T(k) + O(n).
 \end{aligned}$$

رابطه‌ی آخر را با استقرا حل می‌کنیم:

فرض: $T(n) \leq cn$

داده ساختارها و مبانی الگوریتم‌ها

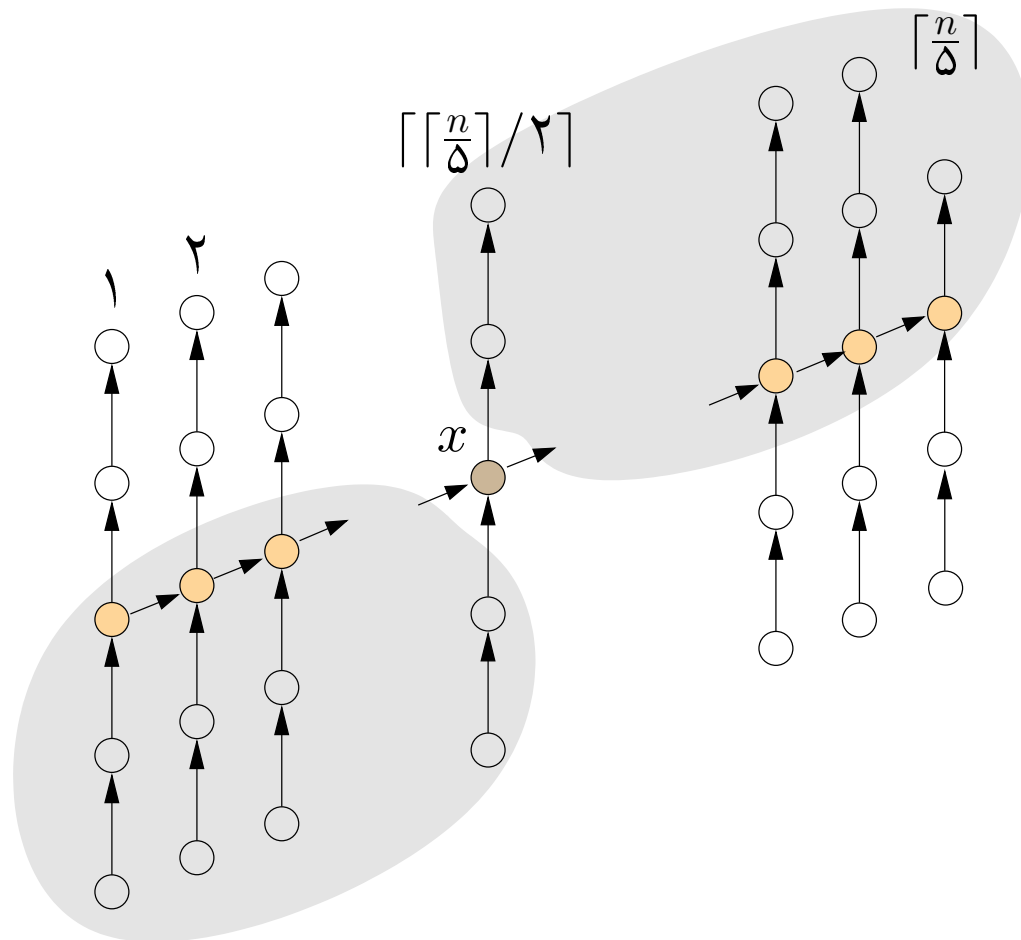
$$\begin{aligned} T(n) &\leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} ck + O(n) \\ &\leq \frac{2c}{n} \left[\sum_{k=1}^{n-1} k - \sum_{k=1}^{\lceil n/2 \rceil - 1} k \right] + O(n) \\ &= \frac{2c}{n} \left[\frac{1}{2}(n-1)n - \frac{1}{2} \left(\left\lceil \frac{n}{2} \right\rceil - 1 \right) \left\lceil \frac{n}{2} \right\rceil \right] + O(n) \\ &\leq c(n-1) - \frac{c}{n} \left(\frac{n}{2} - 1 \right) \left(\frac{n}{2} \right) + O(n) \\ &= c \left[\frac{3}{4}n - \frac{1}{2} \right] + O(n) \\ &\leq cn \end{aligned}$$

زیرا می‌توانیم c را به قدر کافی بزرگ انتخاب کنیم به طوری که $c(n/4 + 1/2) > O(n)$

یافتن k امین عنصر در $O(n)$

محور را طوری انتخاب می‌کنیم تا تضمین کنیم که اندازه‌ی دوبخش در بدترین حالت $O(n)$ هستند.

داده ساختارها و مبانی الگوریتم‌ها



- حداقل $2 - \lceil \frac{n/5}{2} \rceil$ گروه دارای ۳ عنصر کوچک‌تر از x هستند.
- تعداد عناصر کوچک‌تر از x حداقل

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

است.

- به صورت مشابه، تعداد عناصر بزرگ‌تر از x نیز حداقل $3n/10 - 6$ خواهد بود.
- پس، در بدترین حالت، الگوریتم به صورت بازگشتی بر روی حداکثر $7n/10 + 6$ عنصر اعمال خواهد شد.

داده ساختارها و مبانی الگوریتم‌ها

$$T(n) \leq \begin{cases} \Theta(1) & n \leq \Lambda_0 \\ T(\lceil n/5 \rceil) + T(\sqrt{n}/10 + 6) + O(n) & n > \Lambda_0 \end{cases}$$

فرض: برای یک c مفروض و هر $n \leq \Lambda_0$ داشته باشیم $T(n) \leq cn$.

$$\begin{aligned} T(n) &\leq c\lceil n/5 \rceil + c(\sqrt{n}/10 + 6) + O(n) \\ &\leq cn/5 + c + \sqrt{cn}/10 + 6c + O(n) \\ &\leq 9cn/10 + 7c + O(n) \\ &\leq cn \end{aligned}$$