

به نام خدا

الگوریتم Timsort

مهدی حقوردی



فهرست مطالب

مقدمه و معرفی

تاریخچه

چرا تیم سورت؟

توضیح الگوریتم - Run ها

توضیح الگوریتم - ادغام ها

ادغام پیشرفته تر با Galloping

استفاده از stack برای ادغام

روند کدها

مراجع

مقدمه و معرفی

- در دنیای علوم کامپیوتر، مرتب‌سازی یک عملیات اساسی با کاربردهای بی‌شمار است.
- در میان انبوهی از الگوریتم‌های مرتب‌سازی، یکی از الگوریتم‌ها به دلیل کارایی، تطبیق‌پذیری و طراحی زیبا متمایز شده است: الگوریتم تیم‌سورت¹.
- این الگوریتم که توسط تیم پیترز² برای زبان برنامه نویسی پایتون³ توسعه یافته است، به سنگ بنای پیاده‌سازی مرتب‌سازی در زبان‌ها و محیط‌های مختلف برنامه‌نویسی تبدیل شده است.
- ترکیب منحصر به فرد مرتب‌سازی ادغامی⁴ و مرتب‌سازی درجی⁵ به همراه بهینه‌سازی‌های مخصوص روی هر الگوریتم و بهینه‌سازی‌های تطبیقی، تیم‌سورت را به یکی از پیچیده‌ترین و کاربردی‌ترین الگوریتم‌های مرتب‌سازی موجود تبدیل کرده است.

¹ Timsort

² Tim Peters

³ Python programming language

⁴ Merge sort

⁵ Insertion sort

تاریخچه

- الگوریتم تیم‌سورت، در سال ۲۰۰۲ توسعه یافت.

- تیم پیترز این الگوریتم را اینگونه توصیف می‌کند:

“A non-recursive adaptive stable natural mergesort / binary insertion sort hybrid algorithm”

- این الگوریتم از Python 2.3 تا حدود بیست سال، الگوریتم استاندارد مرتب‌سازی در پایتون بود و از نسخه‌ی 3.11.1 به دلیل تغییراتی که در سیاست‌های ادغام آن بوجود آمد، الگوریتمی به اسم Powersort بر پایه‌ی تیم‌سورت، جایگزین آن شد.

- الگوریتم تیم‌سورت در 7 Java SE، Android، GNU Octave، V8، Swift و Rust پیاده‌سازی شده است.

- چرا non-recursive؟
چون طبق گفته‌ی تیم پیترز: «به طور خلاصه، روتین اصلی یک بار از سمت چپ تا راست، آرایه را طی، Runها² را شناسایی و هوشمندانه آنها را با هم ادغام می‌کند.»
- چرا adaptive؟
چون این الگوریتم با توجه به طول و ترتیب‌های از قبل موجود در آرایه، و همچنین بر اساس اندازه‌ی Runهای پیدا شده، تصمیماتی می‌گیرد تا از الگوریتم بهتری برای آن موقعیت استفاده کند.

² در ادامه مفهوم Run توضیح داده می‌شود.

- چرا stable؟

چون این الگوریتم، ترتیب عناصر یکسان در آرایه‌ی اولیه را حفظ می‌کند. برای مثال اگر لیستی از این اسامی داشته باشیم: [peach, straw, apple, spork] و آنرا بخواهیم بر اساس حرف اول کلمات مرتب کنیم، چنین چیزی می‌گیریم: [apple, peach, straw, spork] اگر دقت کنید در لیست اولیه، straw قبل از spork آمده بود و در لیست مرتب شده هم همین ترتیب حفظ شد. به این نگهداری ترتیب اولیه‌ی عناصر، پایداری الگوریتم مرتب‌سازی می‌گویند.

- چرا hybrid؟

چون این الگوریتم از ترکیب دو الگوریتم merge sort و binary insertion sort برای مرتب‌سازی استفاده می‌کند.

چرا تیم سورت؟

چرا تیم سورت؟

- مقایسه‌ی کلی بین پیچیدگی زمانی الگوریتم‌های Quick sort، Merge sort و Heap sort چنین چیز است [Sabah et al., 2023]:

Algorithm	Best Case	Avg. Case	Worst Case	Space
Quick Sort	$O(n \lg n)$	$O(n \lg n)$	$O(n^2)$	$O(\lg n)$
Tim Sort	$O(n)$	$O(n \lg n)$	$O(n \lg n)$	$O(n)$
Merge Sort	$O(n \lg n)$	$O(n \lg n)$	$O(n \lg n)$	$O(n)$
Heap Sort	$O(n \lg n)$	$O(n \lg n)$	$O(n \lg n)$	$O(1)$

- اما این تحلیل کلی یک سری جزئیات راجع به پیچیدگی زمانی الگوریتم را پنهان می‌کند که آن پیچیدگی یک constant factor اثرگذار در میزان پیچیدگی الگوریتم است. $(c_f \cdot n \lg n)$

چرا تیم سورت؟ (ادامه)

- برای مثال در الگوریتم Quick sort انتخاب مقدار left، right و pivot تاثیرگذار است و در n های کوچک سرعت را پایین می آورد.
- در الگوریتم Merge sort هم ما فضایی به اندازه $n + m$ برای ادغام کردن آرایه ها آن هم به صورت بازگشتی و تعداد زیاد نیاز دارد. همچنین این الگوریتم یک الگوریتم بازگشتی است و درخت بازگشتی و یک system stack برای اجرا نیاز دارد.
- بخاطر جابجایی هایی در الگوریتم Heap sort انجام می شود، Locality of Reference در آن نقض شده و پیشبینی های پردازنده برای کش کردن داده ها را تضعیف می کند.

چرا تیم سورت؟ (ادامه)

پس اگر بتوانیم این constant factor را کاهش دهیم
می توانیم سرعت بیشتری از $O(n \lg n)$ بگیریم.

مرتب سازی درجی دودویی

- پیچیدگی زمانی insertion sort برابر با $O(n^2)$ است و constant factor آن بسیار بسیار پایین است چون اولاً inplace عمل می‌کند (پس نیازی به فضای اضافه ندارد) و ثانیاً فقط بین عناصر آرایه پیمایش انجام می‌دهد (پس Locality of Reference هم در آن بسیار خوب است و پردازنده می‌تواند داده‌ها را کش کند).
- در تحلیل‌های انجام شده روی الگوریتم‌ها، این الگوریتم روی تعداد ورودی ۶۴ و پایین‌تر از الگوریتم‌های دیگر مرتب سازی سریع‌تر عمل می‌کند. [Ali et al., 2021]
- الگوریتم binary insertion sort بجای جستجوی خطی در آرایه (با پیچیدگی $O(n)$) در آن جستجوی دودویی انجام داده و در زمان لوگاریتمی ($O(\lg n)$) مکان صحیح آیتم را پیدا می‌کند (علت استفاده از این الگوریتم در ادامه روشن خواهد شد).

مرتب سازی درجی دودویی (ادامه)

- با تعویض نوع جستجوی این الگوریتم میزان پیچیدگی آن (حالت مورد انتظار و در بدترین حالت) تغییری نکرده و همان $O(n^2)$ باقی می ماند؛ اما در CPython مقایسه ها (بخاطر ماهیت dynamic typed بودن زبان) نسبت به جابجا کردن آبجکت ها بسیار وحشتناک کندتر هستند.
- جابجا کردن آبجکت ها صرفا کپی کردن ۸ بایت pointer است اما مقایسه ها میتوانند بسیار کند باشند (چون ممکن است چند متد در سطح پایتون را صدا بزنند) و حتی در حالات ساده ممکن است بین ۳ یا ۴ تصمیم گرفته بشود:
 ۱. تایپ عمل وند چپ چیست؟
 ۲. تایپ عمل وند راست چیست؟
 ۳. آیا باید آنها را به یک تایپ مشخص تبدیل کرد؟
 ۴. چه کدی برای مقایسه این دو موجود هست؟ ...

مرتب سازی درجی دودویی (ادامه)

- پس یک مقایسه ساده باعث تعداد بسیار زیادی C-level pointer dereference، عملیات‌های شرطی و صدا زده شدن توابع می‌شود.
- پس اگر ما تعداد مقایسه‌ها کمتر کنیم میتوانیم سرعت مرتب سازی را بیشتر کنیم (که با استفاده از binary insertion sort ما تعداد مقایسه‌ها را کم می‌کنیم).
- این مرتب سازی با تابع `binarysort`¹ انجام می‌شود.

¹ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1242>

اگر آرایه را به تکه‌های کوچک تقسیم کنیم (برای مثال ۳۲ تا ۶۴ تایی) و سپس آنها را جدا جدا با مرتب سازی درجی مرتب کنیم و سپس همه را ادغام کنیم، می‌توانیم سرعت مرتب سازی را افزایش دهیم.

۱. چون از مرتب سازی درجی که برای تعداد کم سریع است استفاده کردیم $(c_i(32 \text{ to } 64)^2)$ ،

۲. ادغام دو آرایه مرتب شده در زمان $O(n)$ انجام می‌شود و

۳. با این کار ما توانستیم ۵ سطح از درختی در Merge sort تولید می‌شود را کم کنیم $(c_t.n[\lg n - 5])$

در نهایت چون مقدار پیچیدگی مرتب سازی درجی کوچک است، پیچیدگی تیم‌سورت چنین می‌شود:

$$T(n) = c_t.n[\lg n - 5]$$

در دنیای واقعی و داده‌های واقعی معمولاً آرایه‌ها اصطلاحاً partially sorted هستند.

- این به این معناست که تکه‌هایی از آرایه از قبل مرتب هستند؛ برای مثال در این آرایه: $[5, 4, 1, 2, 3]$ قسمت $[1, 2, 3]$ از قبل مرتب است.

- یا حداقل به صورت صعودی یا نزولی پشت سر هم حضور دارند؛ برای مثال در این آرایه: $[6, 4, 1, 2, 3, 5, 7]$ قسمت $[1, 2, 3, 5, 7]$ به صورت صعودی و قسمت $[6, 4]$ به صورت نزولی مرتب است.

الگوریتم تیم‌سورت این تکه‌های صعودی و یا اکیدا نزولی را در آرایه پیدا می‌کند و آنها را Run می‌نامد و از مرتب بودن اولیه آرایه برای افزایش سرعت استفاده می‌کند.

- اگر از حقیقت قبلی استفاده کنیم و آرایه را به قسمت‌هایی صعودی و یا اکیدا نزولی تقسیم کنیم می‌توانیم درخت Merge sort حتی بیشتر از قبل هم کوتاه کنیم.
- و پیچیدگی را به $c_t \cdot n[\lg n - x]$ تبدیل کنیم. پس هر چقدر تعداد Run ها بیشتر باشد، مقدار x بیشتر و درخت ادغام ما کوتاه‌تر و الگوریتم سریع‌تر می‌شود.

توضیح الگوریتم - Run ها

- تکه‌هایی از آرایه که از قبل دارای ترتیب هستند، در این الگوریتم یک Run نام می‌گیرند.

- Runها میتوانند:

۱. صعودی باشند: $a_0 \leq a_1 \leq a_2 \leq \dots$

۲. اکیدا نزولی باشند: $a_0 > a_1 > a_2 > \dots$

- دلیل اینکه یک Run باید نزولی اکید باشد تا یک Run شناخته شود اینست که الگوریتم تیم‌سورت Runهای نزولی را به صورت در جا، برعکس می‌کند و اگر در یک Run نزولی (و نه اکیدا نزولی) دو عنصر یکسان وجود داشته باشد، ماهیت stable بودن الگوریتم نقض می‌شود. برای مثال این آرایه: $[4, \textcolor{red}{3}, \textcolor{blue}{3}, 1]$ اگر برعکس شود: $[1, \textcolor{blue}{3}, \textcolor{red}{3}, 4]$ ؛ که ترتیب عناصر ۳ عوض شده است. اما اگر اکیدا نزولی باشد دیگر این مشکل وجود نخواهد داشت.

- نکته‌ی دیگر اینست که Run ها حداقل دو آیتم دارند مگر وقتی که آخرین عضو آرایه را برای Run جدید برگزینیم.
- اگر عناصر آرایه رندوم باشند، بعید است که ما Run های طبیعی (یعنی قسمتی از آرایه که از قبل مرتب شده باشد) بلندی را شاهد باشیم. اگر یک Run طبیعی تعداد عناصرش کمتر از minrun باشد (توضیح داده خواهد شد)، الگوریتم با استفاده از binray insertion sort تعداد آنرا به حداقل اندازه‌ی Run می‌رساند.
- دلیلی که می‌توانیم از binary insertion sort استفاده کنیم اینست که هر Run خودش مرتب است، پس می‌توان در آن جستجوی دودویی انجام داد.

Run ها (ادامه)

- الگوریتم برای پیدا کردن Run ها چنین عمل می‌کند: فرض کنید چنین آرایه‌ای داریم:
[8, 12, 9, 17, 15, -1, 22, 11, 10, 7] و حداقل اندازه‌ی Run ها هم ۳ تعیین شده است،
- از چپ به راست حرکت می‌کنیم و [8] را جدا می‌کنیم و سراغ آیتم بعدی می‌رویم و متوجه می‌شویم که یک Run صعودی داریم: [8, 12] ادامه می‌دهیم و به عدد ۹ می‌رسیم، چون این Run باید صعودی باشد و ۹ از ۱۲ کمتر است با استفاده از binary insertion sort جایگاه ۹ را پیدا می‌کنیم: [8, 9, 12]
- عدد بعدی هم بزرگ‌تر از ۱۲ است و آن را هم به این Run اضافه می‌کنیم: [8, 9, 12, 17]. عدد بعدی، از ۱۷ کمتر است و چون ما حداقل اندازه‌ی یک Run را داریم آنرا دیگر به این Run اضافه نمی‌کنیم و به سراغ Run بعدی می‌رویم.
- Run بعدی چنین روندی دارد: [15] سپس [15, -1] و سپس با binary insertion sort: [22, 15, -1] و چون ۱۱ از ۱- بیشتر است و ما حداقل اندازه‌ی یک Run را داریم به سراغ Run بعدی می‌رویم؛ که Run بعدی هم این است: [11, 10, 7]

- اگر داده‌ها رندوم باشند، اکثر Run ها یک اندازه خواهند داشت که دو خوبی دارد:

۱. ادغام کردن Run هایی که اندازه‌ی برابر دارند بسیار بهینه است و

۲. ما حداقل توانسته‌ایم اندازه‌ی درخت بازگشتی ادغام را به اندازه‌ی $\log(\text{minrun})$ کم کنیم.

- برای داده‌های واقعی هم، ما چون Run های نسبتاً بلندی خواهیم داشت توانسته‌ایم کوتاه‌ترین درخت بازگشتی ادغام را داشته باشیم و در نتیجه تعداد ادغام‌ها را کم کنیم.

- پیدا کردن Run ها در آرایه توسط تابع `count_run()`¹ انجام می‌شود.
- قسمتی که به Run های اضافه شده، عنصر اضافه می‌کند تا به اندازه `minrun` برسند:²

```

2410      /* If short, extend to min(minrun, nremaining). */
2411      if (n < minrun) {
2412          const Py_ssize_t force = nremaining <= minrun ?
2413                                  nremaining : minrun;
2414          if (binarysort(&ms, lo, lo.keys + force, lo.keys + n) < 0)
2415              goto fail;
2416          n = force;
2417      }

```

¹ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1316>

² <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L2410>

- اگر طول آرایه کمتر از ۶۴ باشد، از binary insertion sort برای مرتب کردن آن استفاده می‌شود.
- اگر طول آرایه توانی از ۲ بود، طبق تست‌های انجام شده تمامی اعداد ۸ و ۱۶ و ۳۲ و ۶۴ و ۱۲۸ سرعت یکسانی را به الگوریتم می‌دادند اما مثلاً در اندازه‌ی ۲۵۶ تا جابجا کردن عناصر در مرتب سازی هزینه بردار و در اندازه‌ی ۸ تعداد صدا زده شدن توابع هزینه بردار بود. بعد از کمی مطالعه عدد ۳۲ برای minrun انتخاب شد.
- اما بعد از زمان زیادی یک اشکال در انتخاب این عدد پیدا شد، این مثال را ببینید:
 $\text{divmod}(2112, 32) \rightarrow (66, 0)$
 که اگر با این تعداد Run ما ادغام را انجام بدهیم، در پایان باید یک آرایه‌ی ۲۰۴۸ عضوی و یک آرایه‌ی ۶۴ عضوی را با هم ادغام کنیم که اصلاً خوب نیست.
- اما اگر عدد ۳۳ را برای minrun انتخاب کنیم ما ۶۴ تا Run با اندازه‌ی ۳۳ داریم که برای ادغام خیلی بهتر است.

- سیاستی که برای محاسبه‌ی minrun در پیش گرفته شده است اینست که این مقدار از `range(32, 64)` به صورتی انتخاب می‌شود که N/minrun یا دقیقاً توانی از دو باشد، یا اگر این ممکن نبود، اکیدا کمتر از ۲ باشد.
- این انتخاب توسط تابع `merge_compute_minrun()`¹ انجام می‌شود.

¹ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L2012>

توضیح الگوریتم - ادغامها

بهرتر کردن استفاده از فضا در ادغام

- الگوریتم ادغام در Merge sort الگوریتمی inplace نیست و برای ادغام کردن دو آرایه با طولهای n و m یک آرایه‌ی جدید به طول $n + m$ نیاز دارد،
- اما تیم‌سورت برای inplace مرتب کردن و کمتر کردن space overhead الگوریتم مرتب سازی ادغامی از روش بهتری برای ادغام کردن دو Run استفاده می‌کند.

بهرتر کردن استفاده از فضا در ادغام (ادامه)

- اگر آرایه‌ای به این صورت داشته باشیم:

$$A = [\underbrace{12, 19, 21, 22}_X, \underbrace{3, 5, 17, 22, 107, 109}_Y]$$

دو Run به X و Y در آن می‌توان شناسایی کرد.

- در این روش جدید ادغام ما یک آرایه‌ی موقت به اندازه‌ی Run کوچک‌تر درست می‌کنیم و عناصر Run کوچک‌تر را در آن کپی می‌کنیم. کپی کردن‌های این چنین هم (یعنی یک تکه از یک آرایه) به بهینه‌ترین حالت در سطح پایین اجرا می‌شوند. پس حالا ما چنین چیزی داریم:

$$A = [12, 19, 21, 22, 3, 5, 17, 22, 107, 109]$$

$$T = [12, 19, 21, 22]$$

- چون آرایه‌ی کوچک‌تر سمت چپ بود، پس از سمت چپ آرایه‌ی موقت و دومین Run شروع به ادغام کردن می‌کنیم. که یعنی:

$$A = [12, 19, 21, 22, \textcircled{3}, 5, 17, 22, 107, 109]$$

$$T = [\textcircled{12}, 19, 21, 22]$$

را باهم مقایسه، و برنده (عنصر کوچک‌تر) را در ایندکس صفر آرایه‌ی A می‌گذاریم. که می‌شود:

$$A = [3, 19, 21, 22, \textcircled{3}, 5, 17, 22, 107, 109]$$

$$T = [\textcircled{12}, 19, 21, 22]$$

گام‌ها (ادامه)

- سپس:

$$A = [3, 19, 21, 22, 3, \textcircled{5}, 17, 22, 107, 109]$$

$$T = [\textcircled{12}, 19, 21, 22]$$

- که می‌شود:

$$A = [3, 5, 21, 22, 3, \textcircled{5}, 17, 22, 107, 109]$$

$$T = [\textcircled{12}, 19, 21, 22]$$

گام‌ها (ادامه)

- سپس:

$$A = [3, 5, 21, 22, 3, 5, (17), 22, 107, 109]$$

$$T = [(12), 19, 21, 22]$$

- که می‌شود:

$$A = [3, 5, 12, 22, 3, 5, (17), 22, 107, 109]$$

$$T = [(12), 19, 21, 22]$$

گام‌ها (ادامه)

- سپس:

$$A = [3, 5, 12, 22, 3, 5, (17), 22, 107, 109]$$

$$T = [12, (19), 21, 22]$$

- که می‌شود:

$$A = [3, 5, 12, 17, 3, 5, (17), 22, 107, 109]$$

$$T = [12, (19), 21, 22]$$

گام‌ها (ادامه)

- سپس:

$$A = [3, 5, 12, 17, 3, 5, 17, \textcircled{22}, 107, 109]$$

$$T = [12, \textcircled{19}, 21, 22]$$

- که می‌شود:

$$A = [3, 5, 12, 17, 19, 5, \textcircled{22}, 22, 107, 109]$$

$$T = [12, \textcircled{19}, 21, 22]$$

گام‌ها (ادامه)

- سپس:

$$A = [3, 5, 12, 17, 19, 5, 17, \textcircled{22}, 107, 109]$$

$$T = [12, 19, \textcircled{21}, 22]$$

- که می‌شود:

$$A = [3, 5, 12, 17, 19, 21, 17, \textcircled{22}, 107, 109]$$

$$T = [12, 19, \textcircled{21}, 22]$$

گام‌ها (ادامه)

- سپس:

$$A = [3, 5, 12, 17, 19, 21, 17, \textcircled{22}, 107, 109]$$

$$T = [12, 19, 21, \textcircled{22}]$$

- که در اینجا عنصر ۲۲ آرایه‌ی موقت انتخاب می‌شود که برای حفظ پایداری الگوریتم مرتب سازی لازم است که بیست و دویی که زودتر آمده را انتخاب کند:

$$A = [3, 5, 12, 17, 19, 21, 22, 22, 107, 109]$$

$$T = [12, 19, 21, \textcircled{22}]$$

- و حالا که عناصر آرایه‌ی موقت تمام شده، باقی عناصر باقی مانده سر جای خود هستند و دیگر مقایسه و جابجایی لازم نیست و آرایه به صورت مرتب ادغام شد:

$$A = [3, 5, 12, 17, 19, 21, 22, 22, 107, 109]$$

- در مثال قبلی Run کوچک‌تر در سمت چپ بود و ما ادغام را از سمت چپ Run بزرگ‌تر و آرایه‌ی موقت انجام دادیم، اما اگر چنین حالتی باشد:

$$A = [\underbrace{3, 5, 17, 22, 107, 109}_X, \underbrace{12, 19, 21, 22}_Y]$$

- ما چنین چیزی داریم:

$$A = [3, 5, 17, 22, 107, 109, 12, 19, 21, 22]$$
$$T = [12, 19, 21, 22]$$

جهت ادغام (ادامه)

- در این حالت که Run کوچکتر در سمت راست بود، ما هم تمام ادغام را از سمت راست انجام می‌دهیم

$$A = [3, 5, 17, 22, 107, \textcircled{109}, 12, 19, 21, 22]$$

$$T = [12, 19, 21, \textcircled{22}]$$

- که می‌شود:

$$A = [3, 5, 17, 22, 107, \textcolor{blue}{109}, 12, 19, 21, \textcolor{red}{109}]$$

$$T = [12, 19, 21, \textcircled{22}]$$

جهت ادغام (ادامه)

- و سپس:

```
A = [3, 5, 17, 22, 107, 109, 12, 19, 21, 109]
T = [12, 19, 21, 22]
```

- که می‌شود:

```
A = [3, 5, 17, 22, 107, 109, 12, 19, 107, 109]
T = [12, 19, 21, 22]
```

و الی آخر...

- تابعی که عمل ادغام را شروع می‌کند، تابع `merge_collapse()`¹ است که در ادامه بیشتر راجع به آن توضیح داده خواهد شد.

¹ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1956>

ادغام پیشرفته‌تر با Galloping

تکنیک Galloping

- چنین سناریویی را تصور کنید: ما این آرایه را داریم:

$$A = [\underbrace{1, 2, 3, 6, 10}_X, \underbrace{4, 5, 7, 9, 12, 14, 17}_Y]$$

- در تکنیک Galloping، ما با استفاده از binary search مکان اولین عضو Y را در X پیدا می‌کنیم که می‌بینیم در جایگاه ۴ باید آن را بگذاریم، این یعنی تمامی عناصر $[1, 2, 3]$ در جایگاه درستی قرار دارند.

- همین کار را برای آخرین عنصر X نسبت به Y انجام می‌دهیم که مکان آن در جایگاه پنجم دومین Run است، این یعنی تمامی عناصر $12, 14, 17$ در جایگاه درستی قرار دارند.

- پس ما چنین آرایه‌ای داریم:

$$A = [1, 2, 3, \underbrace{6, 10}_X, \underbrace{5, 7, 9}_Y, 12, 14, 17]$$

و صرفاً باید الگوریتم ادغامی که بالاتر بحث شد را روی این X و Y کوچک‌تر انجام بدهیم که هم سریع‌تر است و هم فضای کمتری استفاده خواهد کرد.

تکنیک Galloping (ادامه)

- نکته‌ای که در تکنیک Galloping وجود دارد اینست که این روش فقط بعضی وقت‌ها فقط به نفع ماست. به این مثال دقت کنید:

$$A = [\underbrace{1, 3, 5, 7, 9}_X, \underbrace{2, 4, 6, 8, 10}_Y]$$

- ابتدا مکان ۲ را در X پیدا می‌کنیم که در جایگاه دوم باید قرار بگیرد و سپس مکان ۹ را پیدا می‌کنیم که در جایگاه پنجم باید قرار بگیرد.

- که به چنین آرایه‌ای می‌رسیم:

$$A = [1, \underbrace{3, 5, 7, 9}_X, \underbrace{2, 4, 6, 8}_Y, 10]$$

که مشخصا هیچ بهبودی را برای ما نداشت.

تکنیک Galloping (ادامه)

- طبق تست‌هایی که تیم پیترز انجام داده است، تعداد عنصری که باید سر جای خودشان باشند، در ابتدای الگوریتم، باید حداقل ۷ تا باشد که این عدد نام `min_gallop` دارد. این یعنی اگر جایگاه اولین عنصر `Y` را در `X` پیدا کنیم و ۷ عنصر در جایگاه خودشان باشند، این تکنیک موثر خواهد بود (همین برای پیدا کردن آخرین عنصر `X` در `Y` هم صدق می‌کند).
- در ادغام کردن اگر حالت اول (مکان اولین عنصر `Y` در `X`) یا حالت دوم (مکان آخرین عنصر `X` در `Y`) بزرگ‌تر مساوی `min_gallop` بود، از این تکنیک استفاده می‌شود.
- بعد از هر بار موفق بودن `galloping` عدد `min_gallop` به عنوان یک تشویق برای الگوریتم، یکی کم می‌شود تا در ادامه هم بتوان از `galloping` استفاده کرد.

تکنیک Galloping (ادامه)

- حالت عکس هم چنین است که مکان پیدا شده کمتر از `min_gallop` باشد که یعنی این تکنیک موثر نبوده و از حالت مقایسه‌ی ادغامی ساده، یعنی عنصر به عنصر استفاده می‌شود.
- روش galloping در داده‌های واقعی که `partially sorted` هستند بسیار کارآمد است و
 ۱. تعداد مقایسه‌ها را کم و
 ۲. از کپی کردن بهینه سطح پایین استفاده می‌کند.
- اما در داده‌های رندوم این روش موثر نیست و از مقایسه‌ی عنصر به عنصر استفاده می‌شود که
 ۱. کد ساده‌ای دارد و
 ۲. Locality of Reference خوبی هم دارد.

- gallop کردن از چپ توسط تابع `gallop_left()`¹ و gallop کردن از سمت راست توسط تابع `gallop_right()`² انجام می‌شود.
- این دو تابع انعکاس آئینه‌ای همدیگر هستند.

¹ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1371>

² <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1459>

استفاده از stack برای ادغام

ساختار MergeState

- ساختاری در کدهایی که مربوط به الگوریتم تیمسورت پیدا می‌شود به اسم MergeState؛ این ساختار اطلاعات مورد نیاز برای توابع ادغام را داراست.
- اطلاعات دیگری که در این ساختار نگهداری می‌شوند:
 ۱. مقدار `min_gallop`
 ۲. اندازه‌ی استک Run‌ها: `n`
 ۳. خود استک Run‌ها: `pending`
 ۴. آرایه‌ی موقت ادغام کردن: `temparray`
 ۵. و سه تابع `key_compare`، `key_richcompare` و `tuple_elem_compare` برای سناریوهای مختلف مقایسه کردن بین عناصر آرایه.
- کد این ساختار در اینجا¹ قرار دارد.

¹ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1185>

جلوگیری از اجرای بازگشتی

- توابع بازگشتی وقتی که صدا زده می‌شوند، اطلاعات هر بار صدا زده شدن آنها، در CPython، داخل آبجکت frame آنها (frame، آبجکتی است که اطلاعات و context مورد نیاز برای یک code object که حاوی byte code های پایتون است را فراهم و ذخیره می‌کند) در یک stack ذخیره می‌شوند،
- اما آقای تیم پیترز این الگوریتم را که ماهیت بازگشتی ماندنی دارد را غیر بازگشتی نوشته است و برای اینکار از یک stack مخصوص خودش در سطح کد استفاده می‌کند.
- هر بار که یک Run پیدا می‌شود و اندازه‌ی آن به minrun و یا بیشتر می‌رسد، اطلاعات آن روی stack گذاشته می‌شود.
- اگر بخواهیم ساده نگاه کنیم فرض می‌کنیم که خود Run روی استک گذاشته می‌شود اما در واقع فقط آدرس شروع و اندازه‌ی آن روی استک ذخیره می‌شود تا از تکه تکه کردن آرایه و اشغال فضای بیشتر جلوگیری شود.

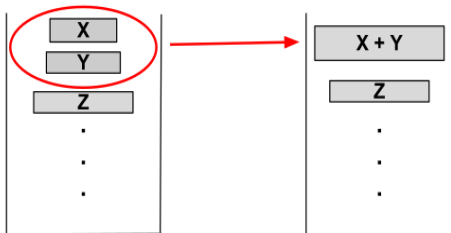
- این استک، نقاط شروع و طول هر Run را در خودش ذخیره می‌کند.
- این یعنی:
Run i ام از آدرس $base[i]$ شروع شده، $len[i]$ تا ادامه دارد و این عبارت همیشه درست است:
 $pending[i].base + pending[i].len == pending[i+1].base$

شروط استک pending

- برای اینکه اندازهی Run ها دقیقا و یا تقریبا با هم برابر بمانند تا الگوریتم های ادغام آنها را سریع تر ادغام کنند، هر بار که یک Run جدیدی ساخته می شود اطلاعات آن روی استک گذاشته می شود، دو شرط برای Run های روی استک چک می شوند.
- اگر روی استک ما Run های X ، Y و Z وجود داشته باشند، این شروط باید حفظ شوند:
 ۱. $|Z| > |Y| + |X|$
 ۲. $|Y| > |X|$
- شرط دوم بر این دلالت دارد که اندازهی Run ها از بالا به پایین به صورت صعودی باشد و
- شرط اول بر این دلالت دارد که وقتی از بالا به پایین می رویم طول Run ها حداقل با سرعت دنبالهی فیبوناچی رشد کنند.
- اگر هر کدام از این شروط نقض بشوند، Run عه Y با Run کوچک تر از بین X یا Z ادغام می شود و شروط دوباره چک می شوند و اگر برقرار بودند، الگوریتم به دنبال Run بعدی در آرایه جستجو می کند.

شروط استک pending (ادامه)

- اگر $|Z| \leq |Y| + |X|$ هر کدام از X یا Z که کوچکتر باشند، با Y ادغام می‌شوند (که الگوریتم X را ترجیح می‌دهد؛ بخاطر اینکه در کش تازگی دارد).



شروط استک pending (ادامه)

- برای مثال اگر این Run ها را داشته باشیم: $\langle X : 10 \quad Y : 20 \quad Z : 30 \rangle$ ، Y با X ادغام می شود:
 $\langle Z : 30 \quad YX : 30 \rangle$
- و یا اگر $\langle X : 1000 \quad Y : 400 \quad Z : 500 \rangle$ ، Y با Z ادغام می شود:
 $\langle C : 1000 \quad ZY : 900 \rangle$
- که در هر دو مثال شرط دوم نقض می شود و الگوریتم تا زمانی که این دو شرط درست شوند، به ادغام کردن ادامه می دهد.
- و در نهایت، تنها یک Run روی استک باقی می ماند که آرایه ی مرتب شده ی ماست.

- تابع `merge_collapse()`¹ مدیریت این شروط را به عهده دارد و انتخاب می‌کند که کدام Run ها با هم ادغام شوند و این کار را با تابع `merge_at()` انجام می‌دهد.
- تابع `merge_at()`² عمل ادغام تصمیمات مبنی بر اینکه ادغام باید با `galloping` شود یا نشود را از نتایج توابع `galloping` اتخاذ و عملیات ادغام با توجه به شرایط را با کمک توابع `merge_lo()`³ و `merge_hi()`⁴ انجام می‌دهد.

¹ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1955>

² <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1891>

³ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1620>

⁴ <https://github.com/python/cpython/blob/3.10/Objects/listobject.c#L1751>

روند کدها

کد کامپایل می‌شود

- فرض کنید شما متد `list.sort()` را در خطی از کدهای خودتان صدا زده‌اید:

```
1 data = list(generate_data())
2 data.sort()
```

- چنین بایت‌کدی برای کد شما زمان کامپایل تولید شده است:

1	1	0	LOAD_NAME	0	(list)
2		2	LOAD_NAME	1	(generate_data)
3		4	CALL_FUNCTION	0	
4		6	CALL_FUNCTION	1	
5		8	STORE_NAME	2	(data)
6		10	LOAD_NAME	2	(data)
7		12	LOAD_METHOD	3	(sort)
8		14	CALL_METHOD	0	
9		16	POP_TOP		
10		18	LOAD_CONST	0	(None)
11		20	RETURN_VALUE		

متد load می‌شود

اول نام متد گرفته می‌شود، و سپس توسط `_PyObject_GetMethod` پیدا می‌شود. اگر نبود error اتفاق می‌افتد.

```
4122     case TARGET(LOAD_METHOD): {
4123         /* Designed to work in tandem with CALL_METHOD. */
4124         PyObject *name = GETITEM(names, oparg);
4125         PyObject *obj = TOP();
4126         PyObject *meth = NULL;
4127
4128         int meth_found = _PyObject_GetMethod(obj, name, &meth);
4129
4130         if (meth == NULL) {
4131             /* Most likely attribute wasn't found. */
4132             goto error;
4133         }
```

چون CPython یک ماشین مجازی stack based است، آبجکت self و خود متد، روی استک push می شوند.

```
4135         if (meth_found) {
4136             /* We can bypass temporary bound method object.
4137              * meth is unbound method and obj is self.
4138              *
4139              * meth | self | arg1 | ... | argN
4140              */
4141             SET_TOP(meth);
4142             PUSH(obj); // self
4143         }
```

```
4159     case TARGET(CALL_METHOD): {  
4160         /* Designed to work in tandem with LOAD_METHOD. */  
4161         PyObject **sp, *res, *meth;  
4162  
4163         sp = stack_pointer;  
4164  
4165         meth = PEEK(oparg + 2);
```



حالا متغیرهای لازم به `call_function` پاس داده می شوند تا کار صدا زدن متد را انجام بدهد.

```
4186      /* This is a method call.  Stack layout:
4187
4188          ... | method | self | arg1 | ... | argN
4189                                     ^- TOP()
4190                               ^- (-oparg)
4191                         ^- (-oparg-1)
4192                   ^- (-oparg-2)
4193
4194      `self` and `method` will be POPed by call_function.
4195      We'll be passing `oparg + 1` to call_function, to
4196      make it accept the `self` as a first argument.
4197      */
4198      res = call_function(tstate, &trace_info, &sp, oparg + 1, NULL);
4199      stack_pointer = sp;
```

تابع call_function این تابع را صدا می زند.

```
2482     int
2483     PyList_Sort(PyObject *v)
2484     {
2485         if (v == NULL || !PyList_Check(v)) {
2486             PyErr_BadInternalCall();
2487             return -1;
2488         }
2489         v = list_sort_impl((PyListObject *)v, NULL, 0);
2490         if (v == NULL)
2491             return -1;
2492         Py_DECREF(v);
2493         return 0;
2494     }
```

مراجع

-  Ali, H., Nawaz, H., and Abdullah, M. (2021).
Performance analysis of heap sort and insertion sort algorithms.
International Journal, 9(5).
-  Sabah, A. S., Abu-Naser, S. S., Helles, Y. E., Abdallatif, R. F., Samra, F. Y. A., Taha, A. H. A., Massa, N. M., and Hamouda, A. A. (2023).
Comparative analysis of the performance of popular sorting algorithms on datasets of different sizes and characteristics.

- ▶ <https://github.com/python/cpython/blob/3.10/Objects/listsort.txt>
- ▶ <https://en.wikipedia.org/wiki/Timsort>
- ▶ http://www.youtube.com/watch?v=emeMe__917E
- ▶ <http://www.youtube.com/watch?v=6D0hQyqAAvU>
- ▶ <http://www.youtube.com/watch?v=Yk4CBisiLaw>
- ▶ <http://www.youtube.com/watch?v=o8y9uYyGLcw>