

به نام خدا

الگوریتم Timsort

مهدی حقوردی



فهرست مطالب

مقدمه و معرفی

تاریخچه

چرا تیمسورت؟

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.

- 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 پایداری الگوریتم مرتب‌سازی می‌گویند.

² در ادامه مفهوم Run توضیح داده می‌شود.

- چرا hybrid?

چون این الگوریتم از ترکیب دو الگوریتم merge sort و binary insertion sort برای مرتب سازی استفاده می کند.

چرا تیم سورت؟

چرا تیم سورت؟

- پیچیدگی زمانی الگوریتم تیم سورت با الگوریتم های Merge sort، Quick sort و Heap sort برابری می کند و برابر $O(n \lg n)$ است.
- اما این تحلیل کلی یک سری جزئیات راجع به پیچیدگی زمانی الگوریتم را پنهان می کند که آن پیچیدگی یک constant factor اثرگذار در میزان پیچیدگی الگوریتم است. $(c_f \cdot n \lg n)$
- 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.

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 ثانياً 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.
- 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

- 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

- 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).

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 $[1, 2, 3, 5, 7]$ خودش به صورت صعودی 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.

- تابع $\text{count_run}()$ ¹ تعداد عنصر موجود در Run را بر می گرداند.

- Run ها 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 اکیدا باشد تا یک Run شناخته شود اینست که الگوریتم تیم سورت Run های نزولی را به صورت در جا، برعکس می کند² و اگر در یک Run نزولی (و نه اکیدا نزولی) دو عنصر یکسان باشند، ماهیت stable بودن الگوریتم نقض می شود. برای مثال این آرایه: $[4, 3, 3, 1]$ اگر برعکس شود: $[1, 3, 3, 4]$ ؛ که ترتیب عناصر ۳ عوض شده است. اما اگر اکیدا نزولی باشد دیگر این مشکل وجود نخواهد داشت.

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#L1064>

- نکته‌ی دیگر اینست که Runها حداقل دو آیتم دارند مگر وقتی که آخرین عضو آرایه را برای Run جدید برگزینیم.
- اگر عناصر آرایه رندوم باشند، بعید است که ما Runهای 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.

- 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 چنین است: $[11, 10, 7]$

- اگر داده‌های رندوم باشند، اکثر Runها یک اندازه خواهند داشت که دو خوبی دارد:

۱. ادغام کردن Runهایی 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 Runهای نسبتاً بلندی خواهیم داشت توانسته‌ایم کوتاه‌ترین درخت بازگشتی ادغام را داشته باشیم و در نتیجه تعداد ادغام‌ها را کم کنیم.

- 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 موقعیت آن بهتر شده 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 است، 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>