

Báo cáo kết quả thực nghiệm các thuật toán sắp xếp

Phan Hoàng Phước

Ngày 11 tháng 3 năm 2023

1 Giới thiệu

Sắp xếp là một trong những giải thuật cơ bản mà bất kỳ lập trình viên nào cũng cần phải biết. Bài báo cáo thực hành IT003 sẽ thực nghiệm một số thuật toán sắp xếp khác nhau và so sánh kết quả chúng để đưa ra cái nhìn tổng quan nhất về mỗi loại thuật toán.

2 Mô tả các thuật toán được đánh giá

2.1 Quick Sort

2.1.1 Khái niệm

Thuật toán nhanh sắp xếp dựa trên việc phân chia mảng dữ liệu thành các nhóm mảng nhỏ hơn, có phương pháp tiếp cận giống [chia để trị](#). Giải thuật chia mảng ra thành 2 phần nhỏ hơn bằng cách so sánh các phần tử đó với một chốt (pivot) là một phần tử có sẵn trong dãy. Một dãy gồm các phần tử lớn hơn chốt và một dãy gồm các phần tử nhỏ hơn chốt.

Thuật toán này sắp xếp dữ liệu nhanh hơn so với nhiều thuật toán sắp xếp truyền thống khác nhưng lại không mấy ổn định vì tùy theo cách chọn khóa (pivot) của người cài đặt.

2.1.2 Ý tưởng

- Thuật toán Quick Sort được thực hiện theo các bước như sau:

1. Nếu độ dài mảng nhỏ hơn hoặc bằng 1 thì kết thúc hàm thực thi hiện tại.
2. Chọn một khóa *pivot* trong đoạn $[l, r]$ cần sắp xếp, gọi vị trí của phần tử này là *id*.
3. Khai báo 2 con trỏ *pleft* và *pright*, *pleft* chạy từ *l* về sau mảng, *pright* từ *r* về trước mảng, lưu ý điều kiện $pleft \leq pright$ luôn thỏa mãn.
4. Thực hiện vòng lặp sau: Tìm vị trí phần tử lớn hơn hoặc bằng vị trí *pleft* lớn hơn *pivot*, tìm vị trí phần tử nhỏ hơn hoặc bằng vị trí *pright* nhỏ hơn *pivot*, hoán đổi giá trị của 2 vị trí đó trong mảng.
5. Sau khi kết thúc vòng lặp ở bước 4, gọi đệ quy 2 dãy từ $[l, id - 1]$ và $[id + 1, r]$ và tiếp tục thực hiện đến hết chương trình.

- Đánh giá:

- Độ phức tạp: nếu ta chọn khóa (*pivot*) một cách ngẫu nhiên thì độ phức tạp trung bình của thuật toán là $\mathcal{O}(n * \log(n))$. Trường hợp tệ nhất của thuật toán là $\mathcal{O}(n^2)$, tuy nhiên ta sẽ không bao giờ gặp trường hợp đó.
- Bộ nhớ tốn: $\mathcal{O}(n)$. **Ưu điểm**
 - Chạy nhanh (nhanh nhất trong các thuật toán sắp xếp dựa trên việc so sánh phần tử).

Nhược điểm

- Phụ thuộc vào cách chọn khóa mà sẽ cho ra thời gian chạy khác nhau (Thông thường sẽ chọn khóa ngẫu nhiên).
- Không ổn định.

2.1.3 Cài đặt

```
1 void QuickSort(float *a, int l, int r) {
2     int pleft = l, pright = r;
3     //chọn khóa random bất kỳ trong đoạn từ [l, r - 1]
4     float pivot = a[l + rand() % (r - l)];
5     //thực hiện vòng lặp
6     while (pleft <= pright) {
7         //tìm vị trí phần tử lớn hơn hoặc bằng vị trí pleft lớn hơn pivot
8         while (a[pleft] < pivot) ++pleft;
9         //tìm vị trí phần tử nhỏ hơn hoặc bằng vị trí pright nhỏ hơn pivot
10        while (a[pright] > pivot) --pright;
11        if (pleft <= pright) {
12            //hoán đổi giá trị tại hai vị trí đó
13            swap(a[pleft], a[pright]);
14            ++pleft;
15            --pright;
16        }
17    }
18    //gọi đệ quy thực hiện sắp xếp 2 dãy con [l, pright] và [pleft, r]
19    if (l < pright) QuickSort(a, l, pright);
20    if (pleft < r) QuickSort(a, pleft, r);
21 }
22
```

2.2 Merge Sort

2.2.1 Khái niệm

Merge Sort là thuật toán sắp xếp trộn theo phương pháp tiếp cận **chia để trị**, bằng cách chia dãy thành 2 dãy con nhỏ hơn có độ dài xấp xỉ nhau và liên tục chia nhỏ chúng ra. Thuật toán giải quyết từ những dãy con nhỏ và sau đó gộp chúng lại và sắp xếp cho dãy lớn hơn.

Merge Sort là một trong thuật toán phổ biến khi lập trình về các bài toán sắp xếp vì tính hiệu quả và dễ cài đặt. Ngoài ra, Merge Sort còn được dùng kết hợp các thuật toán khác để nâng tính hiệu suất trong quá trình sắp xếp.

2.2.2 Ý tưởng

- Thuật toán Merge Sort thực hiện theo các bước như sau:

1. Nếu dãy hiện tại có độ dài là 1 thì trở lại hàm trước.
2. Giả sử dãy hiện tại là đoạn $[l, r]$. Gọi đệ quy 2 dãy con của dãy hiện tại.
3. Sau khi sắp xếp 2 hai dãy con, ta gộp chúng lại vào dãy con hiện tại. Để gộp lại, ta làm như sau:
 - Tạo một dãy mới để lưu các phần sắp xếp vào.
 - So sánh 2 phần tử đầu tiên của 2 dãy con, phần tử nào nhỏ hơn thì ta bỏ vào dãy mới đã được tạo.
 - Tiếp tục thực hiện như vậy đến khi hết dãy.
4. Quay về dãy trước và tiếp tục thực hiện đến khi chạy hết chương trình.

- Đánh giá:

- Độ phức tạp: Do thuật toán chia dãy làm 2 liên tục nên sẽ luôn có độ phức tạp ổn định là $\mathcal{O}(n \cdot \log(n))$ trong mọi trường hợp.
- Bộ nhớ: Do lưu trữ thêm một mảng để lưu lại các phần tử sắp xếp nên sẽ chiếm bộ nhớ là $\mathcal{O}(n)$

Ưu điểm

- Thuật toán chạy nhanh.
- Ổn định với mọi trường hợp.

Nhược điểm

- Tốn thêm 1 mảng lưu trữ.

2.2.3 Cài đặt

```
1 float b[1000001]; //khai báo một mảng lưu trữ
2
3 void MergeSort(float *a, int l, int r) {
4     if (l == r) return;
5     int mid = (l + r) >> 1;
6     //gọi đệ quy 2 mảng con [l, mid] và [mid + 1, r]
7     MergeSort(a, l, mid);
8     MergeSort(a, mid + 1, r);
9     int pl = l, pr = mid + 1, cur = l - 1;
10    while (pl <= mid || pr <= r) {
11        ++cur;
12        //nếu một trong 2 mảng rỗng thì bỏ dãy còn lại vào mảng
13        if (pl > mid) b[cur] = a[pr++];
14        else if (pr > r) b[cur] = a[pl++];
15        else {
16            //thực hiện so sánh 2 phần tử đầu tiên của 2 dãy
17            if (a[pl] > a[pr]) b[cur] = a[pr++];
18            else b[cur] = a[pl++];
19        }
20    }
21    //gán lại vào mảng đang thực hiện sắp xếp
22    for (int i = l; i <= r; i++) a[i] = b[i];
23 }
```

2.3 Heap Sort

2.3.1 Khái niệm

Heap Sort là thuật toán sắp xếp sử dụng cấu trúc dữ liệu Heap. Heap là cấu trúc dữ liệu dựa trên dạng cây nhị phân hoàn chỉnh thỏa mãn tính chất mỗi nút trên cây đều chứa một nhãn lớn hơn hoặc bằng các con của nó (nếu có) và nhỏ hơn hoặc bằng nút cha (trừ nút gốc là và nó là nút lớn nhất). Heap Sort thường được sử dụng nhiều vì có tốc độ nhanh và cài đặt dễ.

2.3.2 Ý tưởng

- Thuật toán Heap Sort hoạt động dựa trên các hoạt động sau đây:

1. Xây dựng một Heap chứa tất cả các phần tử có trong mảng cần sắp xếp.
2. Thực hiện một vòng lặp từ cuối chạy lên đầu và lần lượt gán phần tử tương ứng như sau:
 - Lấy giá trị của nút gốc của Heap ra, đó là giá trị cần gán.
 - Giảm kích thước cây khung đi 1.
 - Nối 2 nhánh cây con (nếu có) sao cho phần tử có giá trị lớn làm gốc.
 - Lặp lại quá trình này đến khi Heap không còn phần tử nào.

- Đánh giá:

- Độ phức tạp: Thuật toán sử dụng cấu trúc dữ liệu Heap nên có chi phí rất lớn và không ổn định. Theo lý thuyết, Heap Sort sẽ có độ phức tạp khi chạy là $\mathcal{O}(n * \log(n))$, tuy nhiên do có chi phí hằng rất lớn nên thời gian chạy có thể chậm hơn.
- Bộ nhớ: $\mathcal{O}(1)$ do cấu trúc dữ liệu Heap được sử dụng dựa trên mảng có sẵn nên không tốn thêm chi phí lưu trữ.

Ưu điểm

- Chạy nhanh.
- Cài đặt đơn giản nếu có sẵn thư viện Heap.

Nhược điểm

- Không ổn định.
- Có chi phí hằng số rất lớn.

2.3.3 Cài đặt

```
1 void heapify(float *a, int n, int i) {
2     // khởi tạo gốc Heap và con trái, con phải của gốc
3     int normalsize = i;
4     int l = 2 * i + 1;
5     int r = 2 * i + 2;
6     // tìm kiếm con lớn hơn nhiều so với gốc
7     if (l < n && a[l] > a[largest]) largest = l;
8     if (r < n && a[r] > a[largest]) largest = r;
9     // nếu giá trị lớn nhất không phải là gốc
10    if (largest != i) {
11        swap(a[i], a[largest]);
12        // đệ quy xây dựng lại Heap
13        heapify(a, n, largest);
14    }
15 }
16
17 void HeapSort(float *a, int l, int r) {
18     int n = r - l + 1;
19     // xây dựng cây heap
20     for (int i = n / 2 - 1; i >= 0; i--) heapify(a, n, i);
21     // sắp xếp lại dãy
22     for (int i = n - 1; i > 0; i--) {
23         // gán phần tử lớn nhất theo vị trí giảm dần
24         swap(a[0], a[i]);
25         heapify(a, i, 0);
26     }
27 }
```

2.4 Hàm sort() của C++

2.4.1 Khái niệm

C++ có thư viện hỗ trợ hàm sort sẵn cho lập trình viên: `std::sort()`. Hàm `sort()` của C++ sử dụng thuật toán Intro Sort, là một thuật toán được kết hợp từ nhiều thuật toán sắp xếp khác để tăng hiệu suất, giảm thời gian chạy như Quick Sort, Heap Sort and Insertion Sort. Đây là thuật toán được sử dụng nhiều nhất vì tính nhanh gọn của hàm C++ được viết sẵn.

2.4.2 Đánh giá

- Độ phức tạp: Thời gian chạy ổn định là $\mathcal{O}(n * \log(n))$
- Bộ nhớ: $\mathcal{O}(1)$

Ưu điểm

- Tốc độ chạy nhanh.
- Dễ dàng cài đặt.

Nhược điểm

- Không ổn định.

2.4.3 Cài đặt

```
1 sort(a, a + n);
```

3 Thực Nghiệm Đánh giá

3.1 Chuẩn bị dữ liệu

Dữ liệu chuẩn bị bao gồm 10 test case khác nhau và code chạy tính thời gian được viết bằng ngôn ngữ C++.

Cụ thể hơn, việc chuẩn bị bao gồm:

1. 10 bộ test case khác nhau gồm 1 triệu số thực dạng float được sinh ngẫu nhiên, trong đó có:
 - Test case 1 là dãy gồm các số thực có giá trị tăng dần.
 - Test case 2 là dãy gồm các số thực có giá trị giảm dần.
 - 8 test case còn lại gồm các thực sinh ngẫu nhiên.
2. 1 đoạn code được viết bằng ngôn ngữ C++ bao gồm các thuật toán
 - Quick Sort
 - Heap Sort
 - Merge Sort
 - `std::sort()`

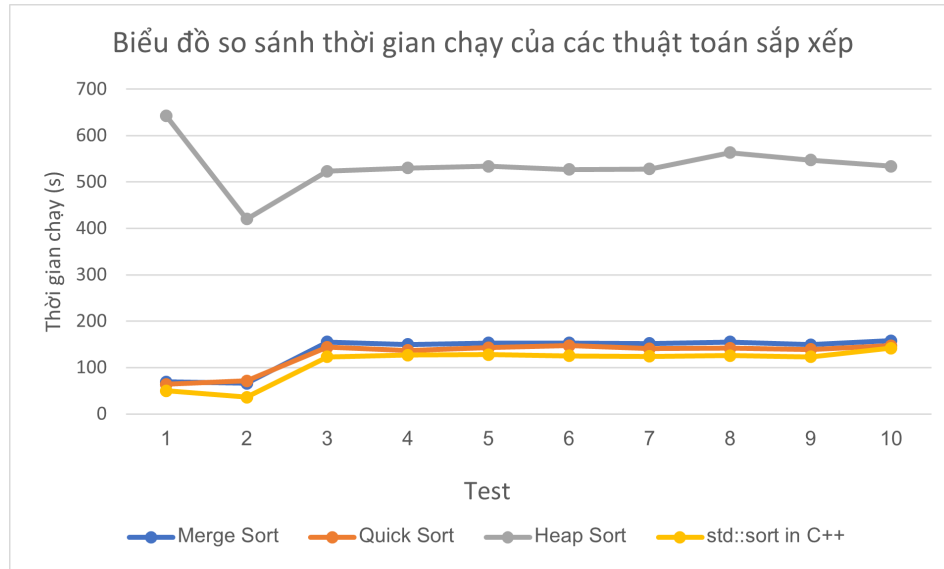
3.2 Kết quả thực nghiệm

Sau khi chạy thực nghiệm các thuật toán, đây là kết quả thống kê được:

Thời gian chạy các thuật toán sắp xếp (ms)				
Test	Merge Sort	Quick Sort	Heap Sort	std::sort in C++
1	69	64	643	50
2	66	71	420	36
3	155	144	523	123
4	150	137	530	127
5	153	143	534	128
6	153	147	527	125
7	152	141	528	124
8	155	142	563	126
9	149	139	547	123
10	158	147	534	142

Hình 1: Bảng dữ liệu thống kê thời gian chạy của các thuật toán sắp xếp

Và khi thực hiện chuyển hóa thành biểu đồ đường so sánh, ta được kết quả như sau:



Hình 2: Biểu đồ đường so sánh thời gian chạy của các thuật toán sắp xếp

3.3 Nhận xét

- Thuật toán Heap Sort có thời gian chạy lâu hơn so với các thuật toán khác do bị tổn chi phí hằng số lớn. Khi mảng sắp xếp tăng dần thì việc xây dựng cấu trúc dữ liệu Heap trở nên tốn thời gian hơn do các phần tử trong Heap liên tục bị đổi gốc. Ngược lại khi dãy sắp xếp giảm dần thì thuật toán chạy hiệu quả và nhanh hơn.
- Các thuật toán khác đều chạy khá ổn định. Đặc biệt là hàm sort của thư viện chuẩn C++ có tốc độ chạy nhanh nhất bởi lẽ đã được tối ưu bằng việc kết hợp nhiều thuật toán khác nhau. Thế nhưng, thời gian chênh lệch giữa các test sinh ngẫu nhiên là khá nhiều vì không ổn định.
- Merge Sort chạy ổn định và không có chênh lệch thời gian nhiều sau mỗi test case được sinh ngẫu nhiên.
- Quick Sort chạy ổn định nhưng vẫn có vài chênh lệch. Bởi lẽ Quick Sort chọn khóa bằng phương pháp chọn ngẫu nhiên trong đoạn nên không ổn định so với Merge Sort.

4 Kết Luận

Thông qua việc thực nghiệm các thuật toán sắp xếp, chúng ta đã hiểu rõ hơn về cách hoạt động, thời gian chạy của các thuật toán theo từng trường hợp khác nhau và có thể rút ra những nhận xét về chúng. Qua đó, tùy theo mục đích và nhu cầu sử dụng mà ta sẽ cài thuật toán phù hợp với chương trình của mình.