

GDSC: Algorithms & Data Structures

Vladislav Artiukhov, based on materials of Sergey Kopelevich

February 27, 2024

Contents

1. Preamble	1
1.1 Credits	2
2. Big-O Notations	3
2.1 Definitions of $O, o, \Omega, \omega, \Theta$	3
2.2 Asymptotic types: linear, quadratic, polylog, exponential	3
3. Basic data structures	4
3.1 Arrays, doubly linked list, singly linked list	4
3.2 <code>std::vector</code> and how it works internally	4
3.3 stack, queue, deque	4
3.4 Keeping minimum for $O(1)$: min-stack, min-queue implementation	4
4. Master Theorem	5
4.1 Master Theorem	6
4.2 Generalized Master Theorem	6
4.3 Algorithm for recurrence relations	6
5. Amortized Analysis	8
5.1 Definition and general perception of the concept	9
5.2 Amortized analysis: Potential method	9
6. Binary Search	11
6.1 Definition, use cases	12
6.1.1 Simplest version	12
6.1.2 Lower bound / Upper bound	12
6.1.3 STL implementation	13
6.2 Use of a predicate	13
6.3 Correctness	14
6.4 Binary search for functions over \mathbb{R}	14

7. Ternary Search	16
7.1 Definition, use cases	17
7.2 Implementation	18
8. Two pointers and Set operations	20
8.1 What a set is	21
8.2 Intersection of sets	21
8.3 Union of sets	21
8.4 Difference of sets	22
8.5 STL implementations	22
8.6 Problem example	22
9. Hash table	23
9.1 Problem statement	24
9.2 Hash table via <code>std::list</code>	24
9.3 Hash functions with probability of collision $O(\frac{1}{m})$	24
9.4 Open addressing collision resolution	25
9.4.1 Hash table overflow	27
9.5 Comparison of the approaches	27
9.6 STL implementation	27

1. Preamble

1.1 Credits

The further work is mostly based on the **Algorithms & Data Structures** course held by Professor **Sergey Kopelevich** in **Higher School of Economics, Saint Petersburg, Applied Mathematics & Computer Science Bachelor's Program**, 1-3 semesters, 2021-2022.

The materials:

1. [SPb HSE, 1st-2nd semesters, Fall 2021/22, Algorithms Lectures Abstract](#)
2. [SPb HSE, 3rd semester, Fall 2021/22, Algorithms Lectures Abstract](#)

Many thanks to Professor **Sergey Kopelevich**!

2. Big-O Notations

2.1 Definitions of $O, o, \Omega, \omega, \Theta$

2.2 Asymptotic types: linear, quadratic, polylog, exponential

For now refer to **GDSC Competitive Programming Abstract (topics 1-5), Fall 2023**.

3. Basic data structures

3.1 Arrays, doubly linked list, singly linked list

3.2 `std::vector` and how it works internally

3.3 stack, queue, deque

3.4 Keeping minimum for $O(1)$: min-stack, min-queue implementation

For now refer to **GDSC Competitive Programming Abstract (topics 1-5), Fall 2023**.

4. Master Theorem

4.1 Master Theorem

Algorithms that are written in a recursive manner oftentimes utilize *divide-and-conquer* technique which implies division of the task into smaller subtasks that are processed by further recursive calls of the algorithm; once the subtask is small enough it is considered as a base case and processed manually. Some examples include **Merge sort algorithm**, **Binary search tree traversal**, etc.

For such algorithms we need to define their asymptotics. **Master Theorem** is a generalized method that yields asymptotically tight bounds for divide and conquer algorithms [wiki].

Theorem 4.1. Master Theorem

Consider the following recurrence relation: $T(n) = a \cdot T(\frac{n}{b}) + f(n)$ where $f(n) = n^c$ for constants $a > 0, b > 1, c \geq 0$; let $k = \log_b n$ be the recursion depth. Then the following holds:

$$\begin{cases} T(n) = \Theta(a^k) = \Theta(n^{\log_b a}), & a > b^c \\ T(n) = \Theta(f(n)) = \Theta(n^c), & a < b^c \\ T(n) = \Theta(k \cdot f(n)) = \Theta(n^c \cdot \log n), & a = b^c \end{cases} \quad (1)$$

Proof.

$$T(n) = f(n) + a \cdot T(\frac{n}{b}) = f(n) + af(\frac{n}{b}) + a^2 f(\frac{n}{b^2}) + \dots + a^k f(\frac{n}{b^k}) \quad | \quad f(n) = n^c$$

$$T(n) = n^c + a \cdot (\frac{n}{b})^c + a^2 \cdot (\frac{n}{b^2})^c + \dots + a^k \cdot (\frac{n}{b^k})^c$$

$$T(b) = n^c \cdot (1 + \frac{a}{b^c} + (\frac{a}{b^c})^2 + \dots + (\frac{a}{b^c})^k)$$

Let $q = \frac{a}{b^c}$ and $S(q) = 1 + q + \dots + q^k$:

1. If $q = 1$: $S(q) = 1 + 1 + \dots + 1 = k + 1 = \log_b n + 1 \implies T(n) = \Theta(f(n) \cdot \log n)$.
2. If $q < 1$: $S(q)$ is a geometric progression, thus it is equal to $S(q) = \frac{1 - q^{k+1}}{1 - q} = \text{const} = \Theta(1) \implies T(n) = \Theta(f(n))$.
3. If $q > 1$: $S(q) = q^k + \frac{q^k - 1}{q - 1} = \Theta(q^k) \implies T(n) = \Theta(a^k \cdot (\frac{n}{b^k})^c) = \Theta(a^k)$.

Note: $f(n)$ could be $O(n^c)$; it does not violate the proof ($f(n) = O(n^c) = C \cdot n^c$).

□

4.2 Generalized Master Theorem

Theorem 4.2. Generalized Master Theorem

In the case of $f(n) = n^c \cdot \log_d n$ Master Theorem still holds:

$$T(n) = a \cdot T(\frac{n}{b}) + n^c \cdot \log_d n, \quad a > 0, \quad b > 1, \quad c \geq 0, \quad d \geq 0.$$

$$\begin{cases} T(n) = \Theta(a^k) = \Theta(n^{\log_b a}), & a > b^c \\ T(n) = \Theta(f(n)) = \Theta(n^c \cdot \log^d n), & a < b^c \\ T(n) = \Theta(k \cdot f(n)) = \Theta(n^c \cdot \log^{d+1} n), & a = b^c \end{cases} \quad (2)$$

4.3 Algorithm for recurrence relations

There are also recurrence relations with the following form:

$$T(n) = a_0 \cdot T(n - p_0) + a_1 \cdot T(n - p_1) + \dots + a_k \cdot T(n - p_k) \quad a_i, p_i > 0, \sum p_i > 1$$

There exists an algorithm of how to find asymptotics for such relations:

Theorem 4.3. Algorithm for recurrence relations

Given $T(n)$ of the above form with the above constants, then the following holds:

$T(n) = \Theta(\alpha^n)$, such that $\alpha > 1$ and it is the **only root** of the equation: $\alpha^n = a_0 \cdot \alpha^{n-p_0} + \dots + a_k \cdot \alpha^{n-p_k}$

Example. Use of Master Theorem

$$1. T(n) = 4 \cdot T\left(\frac{n}{2}\right) + 20 \cdot n^{\frac{3}{2}}$$

$$a = 4, b = 2, c = \frac{3}{2}, f(n) = 20 \cdot n^{\frac{3}{2}} \implies a = 4 > b^c = \sqrt{8} \implies T(n) = \Theta(n^{\log_b a}) = \Theta(n^2).$$

$$2. \text{ Merge sort algorithm recurrence relation: } T(n) = 2 \cdot T\left(\frac{n}{2}\right) + C \cdot n^1$$

$$a = 2, b = 2, c = 1 \implies a = 4 = b^c = 2^1 \implies T(n) = \Theta(n^1 \cdot \log n)$$

Example. Use of Algorithm for recurrence relations

$$1. T(n) = T(n - 1) + 6 \cdot T(n - 2)$$

$$T(n) = \Theta(\alpha^n), \text{ notice that } \alpha = 3 \text{ satisfies the equation: } 3^n = 3^{n-1} + 6 \cdot 3^{n-2}.$$

$$2. T(n) = T(n - 1) + T(n - 2) + T(n - 3)$$

$$1 = \alpha^{-1} + \alpha^{-2} + \alpha^{-3} \implies \alpha \approx 1.839$$

5. Amortized Analysis

5.1 Definition and general perception of the concept

Remember that we were talking about `std::vector` we said that its `push_back` operation works for an average of $\Theta(1)$. It was due to presence of 2 **distinct states**:

1. vector has enough capacity to fit the next pushed element.
2. vector does not have enough capacity and has to make itself twice bigger (i.e. reallocating a memory buffer of size $2 \cdot N$).

There are definitely more complicated scenarios where number of such *interesting states* is much greater, thus we need a unified approach of how to define this average, or **amortized**, time for an operation.

Definition 5.1. The amortized analysis

The amortized analysis is an approach that allows to determine an average running time (time complexity) of operations o_1, o_2, \dots, o_k in a sequence S over that sequence S .

There are several methods that are referred to as amortized analysis ([wiki](#)). We are going to discuss **Potential method**.

5.2 Amortized analysis: Potential method

Definition 5.2. Potential method

Introduce a **potential function** called $\Phi : \mathbf{S} \rightarrow \mathbf{R}_0^+$ where \mathbf{S} is a set of states of the considered data structure and $\mathbf{R}_0^+ = [0, +\infty)$, i.e. the potential function maps states of the data structure to some non-negative values. As an important edge case for the initial state S_{init} : $\Phi(S_{init}) = 0$.

Let o_i be an individual operation within some sequence of operations named Q . Let S_{i-1} be the state of the considered data structure before the execution of the operation o_i and S_i be the state after the execution of o_i . Let Φ be a chosen potential function, then the amortized time for an operation o_i is defined as follows:

$$T_a(o_i) = T_r(o_i) + (\Phi(S_i) - \Phi(S_{i-1})) \text{ where:}$$

1. $T_a(o)$ - amortized time of the operation.
2. $T_r(o)$ - real/actual time spent on the operation.

Theorem 5.1. Potential method yields an upper bound

The amortized time of a sequence of operations always yields an **upper bound** of the the real/actual time for the considered sequence of operations, i.e.:

$\forall O = o_1, o_2, \dots, o_n$ be a sequence of operations, define:

1. $T_a(O) = \sum_{i=1}^n T_a(o_i)$ - amortized time of the sequence O
2. $T_r(O) = \sum_{i=1}^n T_r(o_i)$ - real time of the sequence O

Then: $T_r(O) \leq T_a(O)$

Proof. As definition for $T_a(O)$ suggests:

$$\begin{aligned} T_a(O) &= \sum_{i=1}^n (T_r(o_i) + \Phi(S_i) - \Phi(S_{i-1})) = T_r(O) + \Phi(S_n) - \Phi(S_0) \\ T_a(O) &= T_r(O) + \underbrace{\Phi(S_n)}_{\geq 0} - \underbrace{\Phi(S_0)}_{=0} \geq T_r(O) \end{aligned}$$

Thus, the amortized time provides an accurate upper bound on the actual time of a **sequence of operations**, even though the amortized time for an individual operation may vary widely from its actual time.

□

Note: Generally speaking, Φ could be any function that satisfies the constraints but the idea is to find such Φ that would represent the closest upper bound of the real/actual time for the considered sequence of operations.

Example. Amortized analysis for `std::vector::push_back` operation

1. Analyse `push_back` operation that expands the vector:

Let $\Phi = 2 \cdot s - N$ where the s is the actual size of a vector and N is its capacity. Remember that the expansion occurs once $s == N$.

Let's consider a `push_back` operation that doubles the size of the vector:

1. $T_r = s + 1 = N + 1$ - because we need to copy all the existing elements in a new buffer + place a new element; here $s = N$ because the extension could only happen once $s = N$.

2. $\Phi_0 = 2 \cdot s - N \geq 0$ - value of the potential function before the operation.

3. $\Phi_1 = 2 \cdot (s + 1) - 2N \geq 0$ - value of the potential function after the operation.

Thus, we have:

$$T_a = T_r + \Phi_1 - \Phi_0 = N + 1 + (2s + 2 - 2N - 2s + N) = N + 1 + (2 - N) = N - N + 3 = \Theta(1).$$

2. Analyse `push_back` operation that does not expand the vector:

Notice that in this case $\Delta\Phi = \Phi_1 - \Phi_0 = \text{const}$ and $T_r = \text{const}$, thus $T_a = \text{const} = \Theta(1)$.

For more examples check the [wiki](#) page.

6. Binary Search

6.1 Definition, use cases

6.1.1 Simplest version

Given a sorted array of size n . We need to find an element x in this array for $O(\log n)$:

```

1 | int find(int l, int r, int x) { // [l,r]
2 |     // a is sorted in the ascending order
3 |     while(l <= r) {
4 |         int m = (l + r) / 2;
5 |         if (a[m] == x) return m;
6 |         else if (a[m] < x) l = m + 1;
7 |         else r = m - 1;
8 |     }
9 |     return -1; // i.e., not found
10| }

```

The time complexity is indeed $O(\log n)$ since on each iteration the $|r - l|$ decreases its value in a half ($|r - l| \rightarrow \frac{|r-l|}{2}$).

The problem of such an implementation is that if there are multiple occurrences of x in the array, the algorithm will return **some index** in a range $[l_{ans}, r_{ans})$ where $\forall i = l_{ans}, \dots, r_{ans} - 1 : a[i] == x$, although we would like to get either of both sides of the range, i.e. either l_{ans} or r_{ans} .

6.1.2 Lower bound / Upper bound

The above problem leads to the following implementations of the binary search algorithm that searches for the forementioned indicies:

1. $\min i : a_i \geq x$, i.e. l_{ans} :

```

1 | int lower_bound(int l, int r, int x) {
2 |     while (r - l > 1) {
3 |         int m = (l + r) / 2;
4 |         if (a[m] < x) l = m; // Notice that we keep the invariant: a[l] < x
5 |         else r = m; // => a[r] >= x
6 |     }
7 |     // (a[l] < m) && (a[r] >= m) => r is a minimum index
8 |     return r;
9 | }

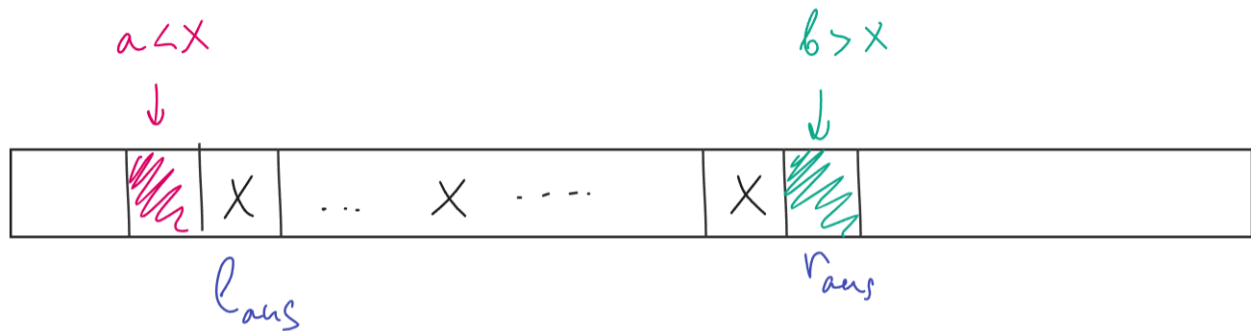
```

2. $\min i : a_i > x$, i.e. r_{ans} :

```

1 | int upper_bound(int l, int r, int x) {
2 |     while (r - l > 1) {
3 |         int m = (l + r) / 2;
4 |         if (a[m] > x) r = m;
5 |         else l = m;
6 |     }
7 |     // Now a[r] > x => l is the last index for which a[l] <= x
8 |     return r;
9 | }

```



6.1.3 STL implementation

In C++ language we are already provided with the template functions that do the same thing:

```

1  /* returns iterator (for int* it is a pointer) */
2  int i = std::lower_bound(a, a + n, x) - a;
3  int i = std::upper_bound(a, a + n, x) - a;
4
5  // For other containers (e.g., std::vector, std::set) that define begin()/end()
   operations use the following:
6  std::lower_bound(
7      std::begin(container),
8      std::end(container),
9      element
10 ) // e.g., for std::vector<T> the return type is std::vector<T>::iterator, i.e.
    pointer to the element

```

6.2 Use of a predicate

We could abstract the implementation even further via searching for a predicate. Predicate is such a function $f : S \rightarrow \{0, 1\}$. Then let's consider the predicate $f(i) = \text{if}(a_i < x) \text{ then } 0 \text{ else } 1$ - in this case the binary search will find such indices l and r that satisfy the following:

1. $l + 1 = r$
2. $f(l) = 0$ (i.e. $a_l < x$)
3. $f(r) = 1$ (i.e. $a_r \geq x$)

```

1  int binary_search(int l, int r, int x) {
2      while(r - l > 1) {
3          int m = (l + r) / 2;
4          if (f(m)) r = m; // invariant: f(r) >= x
5          else l = m; // invariant f(l) < x
6      }
7      return r;
8  }
9
10 // If you want to parameterize predicate as well:
11 // #1: provide a 4th argument as a pointer to a function
12 // (see: https://www.cprogramming.com/tutorial/function-pointers.html)
13 int binary_search(int l, int r, int x, bool(*f)(int)) {
14     ...
15 }
16
17 // #2: template parameter
18 template<typename /* or class */ Func>

```

```

19 | int binary_search(int l, int r, int x, Func f) {
20 |     ...
21 | }
22 | // possible call:
23 | binary_search(l, r, x, []() { /* predicate impl */ });

```

6.3 Correctness

You might already noticed that the functions (arrays are also akin functions, i.e. $a : \{0, 1, \dots, n-1\} : \mathbb{N}$) over which we apply the binary search algorithm are all monotonic functions, i.e. they comply to $x < y : f(x) < f(y)$ (strict monotonically increasing functions; the rest are alike).

Lemma. *The binary search algorithm over some range $[x, y)$ is correct iff the considered function f is monotonic over $[x, y)$.*

6.4 Binary search for functions over \mathbb{R}

It is possible to use the binary search with real numbers as well. Let's consider a problem of *finding square root of x* :

Problem statement: Given $x \in \mathbb{R}$. Find a value $y \in \mathbb{R}$ so that $y^2 == x$ (for us it is $|y^2 - x| < \varepsilon$):

```

1 | double my_sqrt(double x /* never use 'float' */) {
2 |     double l = 0.0;
3 |     double r = x + 1;
4 |
5 |     const double EPS = 1e-9; // 10^(-9)
6 |
7 |     while(r - l > EPS) {
8 |         double m = (r + l) / 2;
9 |         // Preserve invariant: l^2 <= x, r^2 > x
10 |         if (m*m > x) r = m;
11 |         else l = m;
12 |     }
13 |
14 |     return (l + r) / 2;
15 | }

```

In the above notice: if $0 < y < 1$ then $y^2 < y$, and if $y \geq 1$ then $y^2 \geq 1$. Thus, for $0 < x < 1$ the corresponding y is greater than x , i.e. we select $r = x + 1$.

Example. Root of a polynomial $P(x)$

Given a polynomial $P(x)$ of **an odd degree**, i.e. $\deg P = 2k + 1$ with the coefficient for x^{2k+1} be equal to 1. There exists a root $x_0 \in \mathbb{R}$, and we need to find it.

Solution:

We could do that using binary search with any precision ε . First, we need to find points l and r , such that: $P(l) < 0$ and $P(r) > 0$ (e.g., $l = -\infty$, $r = +\infty$ - **MAX_INT** and **MIN_INT** in C++):

```

1 | for (l = -1; P(l) >= 0; l *= 2);
2 | for (r = 1; P(r) <= 0; r *= 2);

```

And finally the root search:

```

1 | while (r - l > EPS) {
2 |     double m = (l + r) / 2;
3 |     if (P(m) < 0) l = m; // P(l) < 0
4 |     else r = m; // P(r) >= 0

```



```
5 ||     }  
6 ||     return (l + r) / 2;
```

Actually we need exactly $k := \frac{r-l}{\varepsilon}$ iterations, thus the while-loop could be changed by the for-loop.

7. Ternary Search

7.1 Definition, use cases

Imagine that we need to find a maximum of a function $f(x)$ in the interval $[l, r]$ but now the function is no longer monotonic. Is there a fast way (i.e. as fast as logarithmic asymptotic as in the binary search algorithm) of finding a point x_0 such that $f(x_0) \rightarrow \max$?

The answer is yes if our function satisfies some constraint.

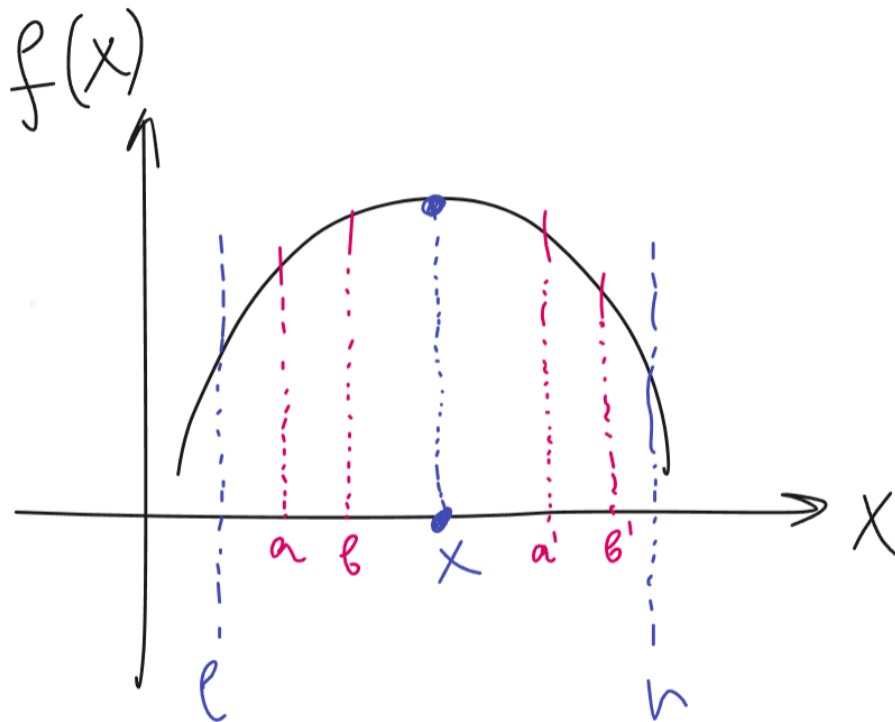
Definition 7.1. Ternary search algorithm

Given a function $f(x)$ and a range $[l, r]$ in which the function is [unimodal](#) and convex, and we need to find maximum in the range (maximum for convex functions, minimum for concave functions).

Since the function is unimodal over $[l, r]$ the following holds:

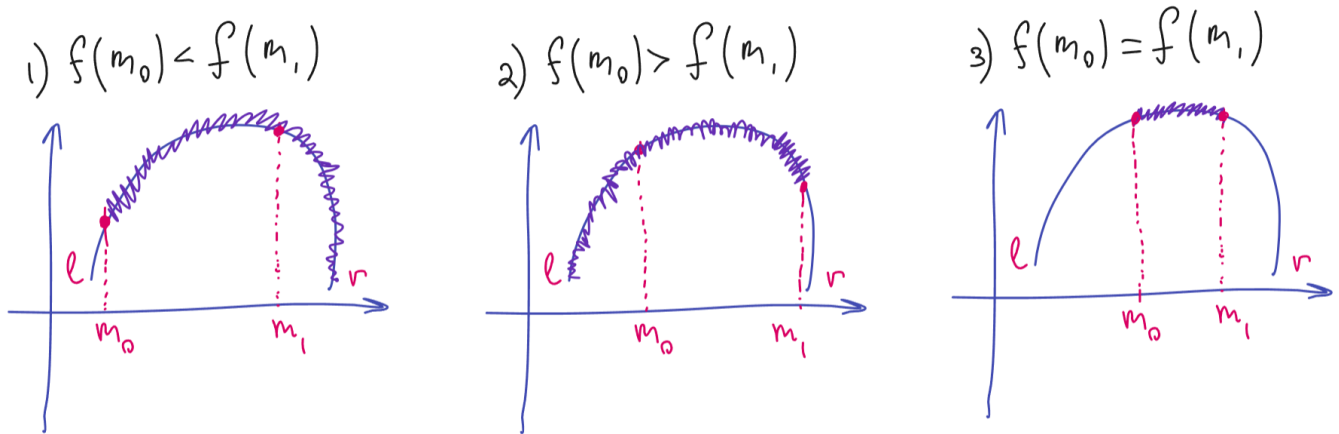
1. $\forall a, b : l \leq a < b \leq x_0 : f(a) < f(b)$
2. $\forall a, b : x_0 \leq a < b \leq r : f(a) > f(b)$

Where x_0 is a point where $f(x) = \max$



For the algorithm let's consider two points m_0 and m_1 such that $l < m_0 < m_1 < r$, there are several cases:

1. $f(m_0) < f(m_1)$, then the maximum of the function f cannot be located in the interval $[l, m_0]$ since there is an interval $[m_0, r]$ on which there are points where function takes greater values.
2. $f(m_0) > f(m_1)$, then the maximum of the function f cannot be located in the interval $[m_1, r]$ since the interval $[l, m_1]$ contains points in which the function takes greater values.
3. $f(m_0) = f(m_1)$, in this case the maximum will be located in the interval $[m_0, m_1]$ but for the implementation simplify this condition branch can be merged with any of the above.



As for the choice of the points m_0 and m_1 we need to select something that would truncate the remaining searching area in a manner that would give us a logarithmic asymptotic. Let's make m_0 be the first 3rd of the interval and m_1 be the second 3rd of the interval:

$$1. m_0 = l + \frac{r-l}{3}$$

$$2. m_1 = l + 2 \cdot \frac{r-l}{3} = \frac{3l+2r-2l}{3} = \frac{2r+l}{3} = \frac{3r-r+l}{3} = r - \frac{r-l}{3}$$

Notice that according to the above statement the searching interval becomes $\frac{2}{3}$ of the initial interval, thus:

$$T(n) = T\left(\frac{2}{3}n\right) + 1$$

According to **Master Theorem**:

$$a = 1, b = \frac{2}{3}, c = 0, f(n) = n^c = 1 \implies a = b^c \implies \Theta(n^c \cdot \log n) = \Theta(\log n) \implies$$

$$T(n) = T\left(\frac{2}{3}n\right) + 1 = \Theta(\log n)$$

7.2 Implementation

Here is the implementation for a function which is an array:

```

1  int ternary_search(vector<int>& arr) {
2      // 'arr' is a unimodal function N -> Z
3      int l = 0;
4      int r = arr.size();
5
6      while(r - l > 2 /* not 1 because m0 = l, m1 = r possible */) {
7          int m0 = l + (r - l) / 3;
8          int m1 = r - (r - l) / 3;
9
10         if (arr[m0] < arr[m1]) {
11             l = m0;
12         }
13         else {
14             r = m1;
15         }
16     }
17
18     int ans = l;
19
20     if (r < arr.size() && arr[r] > arr[ans]) {
21         ans = r;
22     }
23
24     return ans;

```

25 || }

Here is the implementation for a function f over floating-point numbers:

```
1 | int ternary_search(double l, double r, double (*f)(double)) {
2 |     while (r - l > EPS) {
3 |         double m0 = l + (r - l) / 3;
4 |         double m1 = r - (r - l) / 3;
5 |
6 |         if (f(m0) < f(m1)) {
7 |             l = m0;
8 |         }
9 |         else {
10 |             r = m1;
11 |         }
12 |     }
13 |
14 |     return (l + r) / 2;
15 | }
```

8. Two pointers and Set operations

8.1 What a set is

We are going to understand a set as a data structure which contains **distinct** elements (otherwise **multiset**) sorted in the **ascending order**.

STL library contains a data structure that complies to the above definition: `std::set<T>`. In the following subsections we are going to look at fundamental operations over sets. We assume that every time when we are given a set, it is an ascendingly sorted array of distinct elements.

8.2 Intersection of sets

Given two sets A and B and we want to create a set $C = A \cap B$ in $O(|A| + |B|)$:

```

1 vector<int> intersection(vector<int> A, vector<int> B) {
2     int i = 0;
3     int j = 0;
4     vector<int> C;
5
6     for (; i < |A| && j < |B|; ++j) {
7         while (i < |A| && A[i] < B[j]) ++i;
8         // last condition means "either 1st element of B or a distinct element"
9         // condition can be removed since all elements of both A and B are
10        distinct
11        if (i < |A| && A[i] == B[j] && (j == 0 || B[j - 1] != B[j])) {
12            C.push_back(B[j]);
13        }
14    }
15    return C;
16 }
```

The above algorithm can be extended to an intersection of an arbitrary number of sets A_0, A_1, \dots, A_k . You can try to do that in the following LeetCode problem: [2248. Intersection of Multiple Arrays](#).

8.3 Union of sets

Given two sets A and B , we want to find $C = A \cup B$ in $O(|A| + |B|)$:

```

1 vector<int> union(vector<int> A, vector<int> B) {
2     int i = 0;
3     int j = 0;
4     vector<int> C;
5
6     for (; j < B.size(); ++j) {
7         while (i < A.size() && A[i] <= B[j]) {
8             C.push_back(A[i++]);
9         }
10        // in the case if last added element A[i] is equal to B[j]
11        if (C.empty() || C.back() != B[j]) {
12            C.push_back(B[j]);
13        }
14    }
15
16    while (i < A.size()) {
17        C.push_back(A[i++]);
18    }
19
20    return C;
21 }
```

8.4 Difference of sets

Given two sets A and B , we want to find $C = A \setminus B$ in $O(|A| + |B|)$:

```

1 | vector<int> difference(vector<int> A, vector<int> B) {
2 |     int i = 0;
3 |     int j = 0;
4 |     vector<int> C;
5 |
6 |     for (; i < A.size(); ++i) {
7 |         // skipping elements less than A[i]
8 |         while(j < B.size() && B[j] < A[i]) {
9 |             ++j;
10 |        }
11 |        // if A[i] not present in B add A[i] to the answer
12 |        if (j >= B.size() || B[j] != A[i]) {
13 |            C.push_back(A[i]);
14 |        }
15 |    }
16 |
17 |    return C;
18 | }
```

You may solve this task by yourself of LeetCode: [2215. Find the Difference of Two Arrays](#)

8.5 STL implementations

STL already has implementations of all the operations mentioned above: `std::set_difference`, `std::set_union`, `std::set_intersection`, and `std::merge` (the latter one is the same as `std::set_union` but it does not remove the duplicate elements). All the mentioned STL functions are called in the following manner:

```

1 | int k = std::merge(A, A + |A|, B, B + |B|, C) - C;
2 | // C - pointer to where to store the result (memory allocation is your
3 | responsibility)
4 | // k - number of elements in the result, i.e. k == |C|
```

8.6 Problem example

9. Hash table

9.1 Problem statement

We are given an empty set of integers, and we want to be able to quickly add, remove, find operations.

Approach 1: The easiest (and slow) implementation could be to store all the elements in `std::vector`, then:

1. `add = push_back = $O(1)$.`
2. `find = $O(n)$.`
3. `remove = find + $O(1)$ (swap the found element with the last one and perform pop_back).`

Approach 2 (with additional constraint): if all the elements are in the range of $[0, N)$, then we can allocate an array `is[N]`; then `is[x]` indicates whether the element x is present in the set or not, in this case all the required operations will work in $O(1)$.

Idea solution: the hash table is a data structure that allows to perform add, remove, and find in randomized $O(1)$.

9.2 Hash table via `std::list`

```

1 | std::list<int> table[N];
2 |
3 | void add(int x) {
4 |     table[x % N].push_back(x); // O(1)
5 | }
6 |
7 | auto /* = std::list<int>::iterator */ find(int x) {
8 |     const auto& list = table[x % N];
9 |     return std::find(list.begin(), list.end(), x); // works for O(list.size())
10 | }
11 |
12 | void remove(int x) {
13 |     table[x % N].erase(find(x)); // works for find + O(1)
14 | }
```

Instead of `std::list` you can use any container-like data structure, e.g., `std::vector`.

If there are n elements present in the set and all of them are uniformly distributed, then each list will have on average $\frac{n}{N}$ elements \implies if $n \leq N$ + elements are uniformly distributed, then all the operations will work in $O(1)$.

The uniform distribution of the elements is something we can potentially adjust via selecting a good hash function $f(x)$ for the elements (e.g., in the above we have used $f(x) = x \bmod N$).

9.3 Hash functions with probability of collision $O(\frac{1}{m})$

statement 9.1. Functions $f(x) = x \bmod N$ where N is a prime number are good enough.

In the following theorem we will show that there exist a special type of functions whose probability of collision is $O(\frac{1}{m})$ where m is the capacity of the underlying hash table.

Definition 9.1. Universal family of hash functions

Given a universe of elements U .

Let a hash function be defined as: $h : U \rightarrow \{0, 1, 2, \dots, m-1\}$ for some m , or $h : U \rightarrow [m]$ ($[m] := 0, 1, 2, \dots, m-1$).

$H = \{h : U \rightarrow [m]\}$ is called a **universal family** if:

$\forall x, y \in U : x \neq y : |S| \leq \frac{|H|}{m}$ where $S := \{h \in H : h(x) = h(y)\}$. Perceive the $|S|$ as the number of hash functions that collide x and y , then the inequation can be rewritten as:

$\forall x, y \in U : x \neq y : \frac{|S|}{|H|} \leq \frac{1}{m}$ where $\frac{|S|}{|H|}$ is the probability of collision of x and y for the universal family H .

"In other words, any two different keys of the universe collide with probability at most $\frac{1}{m}$ when the hash function h is drawn uniformly at random from H ." (wiki).

Theorem 9.2. Integer hashing

Let's consider a case where the values from the universe U are finite (e.g., integer numbers in 64-bit computers), i.e. $U := \{0, 1, 2, \dots, |U| - 1\}$.

Select a prime p which is $p \geq |U|$ and define:

$h_{a,b}(x) := ((ax + b) \bmod p) \bmod m$ where a, b are randomly chosen integers modulo p , and $a \neq 0$.

The set $H := \{h_{a,b} : a \neq 0, b\}$ is a **universal family**, i.e.:

$\forall x, y \in U : x \neq y : Pr[\text{collision of } x \text{ and } y \text{ occurred}] \leq \frac{1}{m}$.

Proof. Notice that the collision occurs once $h_{a,b}(x) = h_{a,b}(y)$, $x \neq y$ which holds only when:

$$(ax + b) \% p \equiv_m (ay + b) \% p$$

$ax + b \equiv_p ay + b + im$ where $im \leq p - 1$ because the right part of the equation is in $\{0, 1, 2, \dots, p - 1\}$, thus $i = 1, 2, \dots, \frac{p-1}{m}$.

$a(x - y) \equiv_p im$ since $x \neq y$, $x, y \in U = \{0, 1, \dots, |U| - 1\}$ and $p \geq |U|$, and p - prime, then $(x - y) \not\equiv_p 0$, thus the $(x - y)^{-1}$ exists:

$a \equiv_p im \cdot (x - y)^{-1}$ - collision happens once the right and the left parts are equal. For every $i = 1, 2, \dots, \frac{p-1}{m}$ the equality takes place in exactly single value of a because $a = 1, 2, \dots, p - 1$, thus the probability of the collision is $Pr = \frac{\#i}{\#a} = \frac{\frac{p-1}{m}}{p-1} = \frac{1}{m}$.

Thus, we have $Pr[\text{collision of } x \text{ and } y \text{ occurred}] \leq \frac{1}{m} \implies H$ is a universal family. □

9.4 Open addressing collision resolution

The implementation utilizes the cyclic array. A hash function is used to get the initial position in the array of the element x , next we move strictly right (modulo the capacity of the cyclic array) until we find a position in which either no element present or the element x locates:

```

1 | int h[N]; // hash table
2 |
3 | // h[i] = 0 : an empty cell
4 | // h[i] = -1 : deleted element
5 | // h[i] > 0 : there is an element
6 |
7 | int getIndex(int x) { // get index by an element, require x > 0
8 |     int i = x % N;
9 |     while (h[i] && h[i] != x) {
10 |         i = (i + 1) % N;
11 |     }
12 |     return i;
13 | }
14 |
15 | void add(int x) {

```

```

16 |     h[getIndex(x)] = x;
17 | }
18 |
19 | void remove(int x) {
20 |     // notice that the cell does not become empty!
21 |     // marking a cell as empty will break the algorithm! see below
22 |     h[getIndex(x)] = -1; // requires x != -1
23 | }
24 |
25 | int find(int x) {
26 |     return h[getIndex(x)] != x;
27 | }

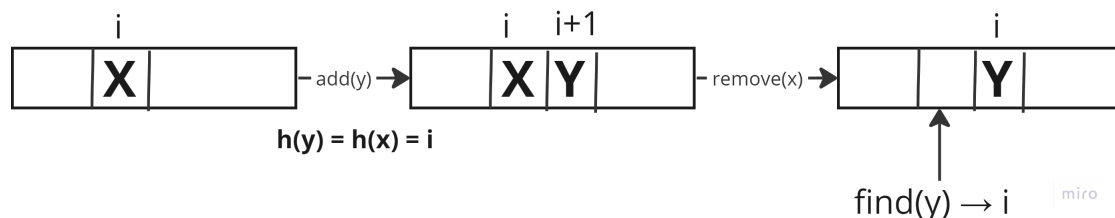
```

Why marking cells as removed is important:

Suppose we mark a cell as empty once element is removed, then let's consider the following scenario:

You have 2 elements x and y which are about to be added into the hash table T , and it turns out that the hash function h over the elements results in $h(x) = h(y)$. Let's execute the following sequence of operations:

1. Suppose T is empty.
2. $add(x)$ - adds x at the position i .
3. $add(y)$ - adds y at the position $i + 1$ because $h(x) = h(y)$, thus the algorithm will traverse right until a free cell is found.
4. $remove(x)$ - marks the position i as empty.
5. $find(y)$ - returns **false** because $h(y) = h(x) = i$, the cell $T[i]$ is marked to be empty, thus the find method says no y found in the T .



Lemma. If the hash table with open addressing of size N has αN occupied cells ($\alpha < 1$), then the mathematical expectation E of the time of **getIndex** $\leq \frac{1}{1-\alpha}$.

For example, if there are only half of the elements occupied, i.e. $\alpha = \frac{1}{2}$, then on average each operation will require $\leq 2 = O(1)$ of time.

Proof. Suppose that the free cells are distributed uniformly due to the hash function being good enough (no proof for it; oftentimes holds if you do not encounter an anti-test for your solution where all elements would have the same hash value).

The worst case is when the element x is not present in the table. In this case, on each iteration of the **while**-loop of **getIndex** the probability to **not stop** equals to $\frac{\alpha \cdot N}{N} = \alpha$ (due to the supposed uniform distribution).

The probability to not stop even after k iterations is α^k , i.e. we will make k -th step with the probability of α^k . Then, the execution time will be:

$$T = \text{mathematical expectation } E \text{ of the number of steps made} = 1 + \sum_{k=1}^{+\infty} Pr[\text{the } k\text{-th step was made}] = 1 + \alpha + \alpha^2 + \alpha^3 + \dots = \frac{1}{1-\alpha}$$

□

9.4.1 Hash table overflow

If α becomes too big the operations start to work slower:

1. $\alpha = 1$: no free cells, **getIndex** will infinitely search for a free cell. How to resolve this problem?

Once $\alpha > \frac{2}{3}$ let's double the size of the underlying cyclic array and in linear time insert the elements into the new allocated buffer. During copying we will of course skip the cells marked as removed (i.e. $T[i] = -1$) \implies the removed cells occupy the additional space only until the next sizing doubling takes place (i.e. there will be no removed cells bloating inside the hash table).

9.5 Comparison of the approaches

Overall, we have discussed two types of the hash table, let's compare them in terms of the required memory and time:

Suppose we store x bytes for an object (element) and a pointer requires 8 bytes, then our hash tables use for n elements:

1. List-based approach (exactly n lists): $8n(\text{pointers to } n \text{ lists}) + n \cdot (x+8)$ (every list has an element and a pointer) $n \cdot (x+16)$ bytes.
2. Open addressing approach (having an additional capacity of $\frac{3}{2}$): $\frac{3}{2} \cdot n \cdot x$ bytes.

Time complexity:

1. List-based solution: requires to make on average 2 lookups, one for a list and one for an element inside the selected list.
2. Open addressing solution: requires on average 1 lookup.

9.6 STL implementation

STL implementation uses list-based approach, thus if you write your own solution with open addressing it will work almost twice faster on average.

1. **std::unordered_set<int> h** - hash table that stores the set of integers.

Use:

1. **std::unordered_set<int> h(N)** - create a hash table with the capacity of N cells.
2. **h.count(x)** - check whether x present.
3. **h.insert(x)** - add an element x (no-op if already present).
4. **h.erase(x)** - remove an element x (no-op if already not present).
2. **std::unordered_map<int, int> h** - hash table that stores **std::pair<int, int>**.

Use:

1. **std::unordered_map<int, int> h(N)** - create a hash table with the capacity of N cells.
2. **h[i] = x** - i -th cell can be used as an array index, stores an association $i \rightarrow x$.
3. **h.count(x)** - check whether a pair with the first component being equal i present (i is a key).
4. **h.erase(x)** - remove a pair with the first component being equal i if present (i is a key).

You may use **std::unordered_map** as an array with arbitrary indices, it is oftentimes called an **associative array** because for each key i there is an association with a value x .

Depending on the set initial capacity of the **std::unordered_map** there exists an anti-test which makes the structure to work in linear time for an operation, thus, if there is a change that a testing

system might have an anti-tests for the initial capacity N , it may be helpful to set the initial capacity to be an random value:

```
1 | int64_t randomTime() {
2 |     std::random_device rd; // Obtain a random number from hardware
3 |     std::mt19937 eng(rd()); // Seed the generator
4 |
5 |     // Define the range of time (in seconds)
6 |     std::uniform_int_distribution<> distr(0, 86400); // 86400 seconds in a day
7 |
8 |     // Get the current time
9 |     std::time_t currentTime = std::time(nullptr);
10 |
11 |     // Generate a random offset from the current time
12 |     std::time_t randomTime = currentTime + distr(eng);
13 |     return static_cast<int64_t>(randomTime);
14 | }
15 |
16 | unordered_map<int, int> h(N + randomTime() % N);
```