

Алгоритмы и структуры данных

Сортировка

д.т.н., проф. Трифонов Петр Владимирович

Содержание лекции

- 1 Постановка задачи
- 2 Сортировка с помощью кучи (пирамиды)
- 3 Приоритетные очереди
- 4 Быстрая сортировка
- 5 Нижняя граница сложности сортировки
- 6 Сортировка подсчетом
- 7 Поразрядная сортировка
- 8 Карманная сортировка
- 9 Медианы и порядковые статистики

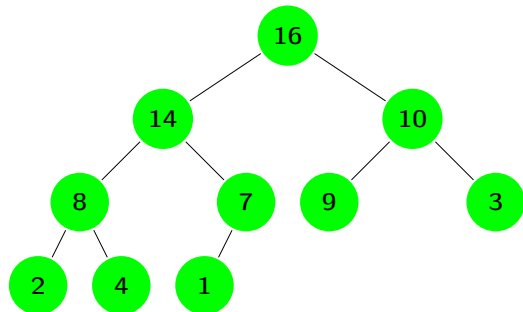
Задача сортировки

- Сортируемые объекты — структуры (*key*, *value*)
- На множестве ключей задано отношение строгого слабого порядка \prec
 - Иррефлексивное: $\forall x : x \not\prec x$
 - Асимметричное: $\forall a, b : a \prec b \Rightarrow b \not\prec a$
 - Транзитивное: $\forall a, b, c : a \prec b \wedge b \prec c \Rightarrow a \prec c$
 - $\forall x, y, z : ((x \not\prec y \wedge y \not\prec x) \wedge (y \not\prec z \wedge z \not\prec y)) \Rightarrow (x \not\prec z \wedge z \not\prec x)$
- Для простоты под $A \prec B$ будем понимать $A.key \prec B.key$
- Сложность измеряется в числе операций сравнения и обмена
- Внутренняя сортировка: весь сортируемый массив загружен в оперативную память
- Внешняя сортировка: данные записаны на медленном носителе, одновременно в памяти находится только часть сортируемых данных
- Устойчивая сортировка: если в исходном массиве $A.key = B.key$ и A расположено раньше B , то в отсортированном массиве A также расположено раньше B

Куча (пирамида)

Двоичная куча — массив (A_0, \dots, A_{n-1}) , для которого $A[\text{Parent}(i)] \geq A[i], 0 < i < n$

- $\text{Parent}(i) = \lfloor \frac{i-1}{2} \rfloor$
- $\text{Left}(i) = 2i + 1$
- $\text{Right}(i) = 2i + 2$
- $A[0]$ — корень дерева
- Высота дерева — число ребер в кратчайшем пути от корня до листа
- Высота вершины поддерева — высота поддерева с корнем в этой вершине
- Высота всего дерева — $\lfloor \log_2(n) \rfloor$



$A = (16, 14, 10, 8, 7, 9, 3, 2, 4, 1)$

Обеспечение основного свойства кучи ($A[\text{Parent}(i)] \succeq A[i], 0 < i < n$)

- Будем считать, что поддеревья с корнями $\text{Left}(i)$ и $\text{Right}(i)$ уже обладают свойством кучи
- Как переставить элементы поддерева с вершиной i , чтобы получилась куча?

Algorithm 1: Heapify(A,i)

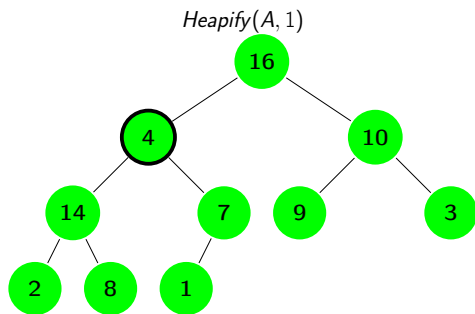
```

1  l=Left(i);r=right(i);
2  if  $l < \text{HeapSize} \wedge A[l] \succ A[i]$  then
3    | largest=l
4  else
5    | largest=i
6  if  $r < \text{HeapSize} \wedge A[r] \succ A[\text{largest}]$  then
7    | largest=r
8  if  $\text{largest} \neq i$  then
9    | swap(A[i],A[largest]);
10   | Heapify(A,largest)

```

Обеспечение основного свойства кучи ($A[\text{Parent}(i)] \succeq A[i], 0 < i < n$)

- Будем считать, что поддеревья с корнями $\text{Left}(i)$ и $\text{Right}(i)$ уже обладают свойством кучи
- Как переставить элементы поддерева с вершиной i , чтобы получилась куча?



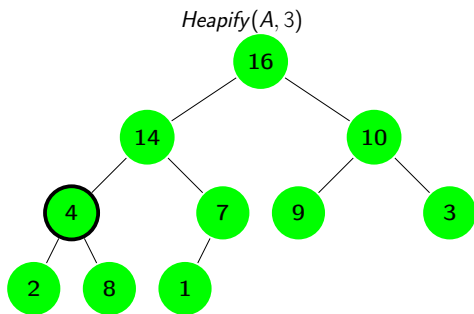
Algorithm 2: Heapify(A,i)

```

1 l=Left(i);r=right(i);
2 if l < HeapSize ∧ A[l] > A[i] then
3   | largest=l
4 else
5   | largest=i
6 if r < HeapSize ∧ A[r] > A[largest] then
7   | largest=r
8 if largest ≠ i then
9   | swap(A[i],A[largest]);
10  | Heapify(A,largest)
  
```

Обеспечение основного свойства кучи ($A[\text{Parent}(i)] \succeq A[i], 0 < i < n$)

- Будем считать, что поддеревья с корнями $\text{Left}(i)$ и $\text{Right}(i)$ уже обладают свойством кучи
- Как переставить элементы поддерева с вершиной i , чтобы получилась куча?



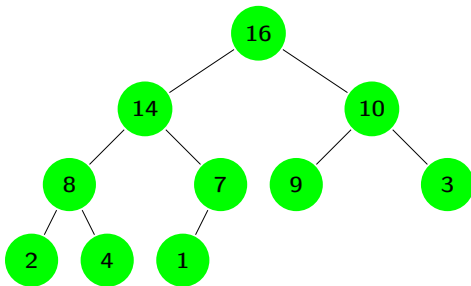
Algorithm 3: Heapify(A,i)

```

1 l=Left(i);r=right(i);
2 if l < HeapSize ∧ A[l] > A[i] then
3   | largest=l
4 else
5   | largest=i
6 if r < HeapSize ∧ A[r] > A[largest] then
7   | largest=r
8 if largest ≠ i then
9   | swap(A[i],A[largest]);
10  | Heapify(A,largest)
  
```

Обеспечение основного свойства кучи ($A[\text{Parent}(i)] \succeq A[i], 0 < i < n$)

- Будем считать, что поддеревья с корнями $\text{Left}(i)$ и $\text{Right}(i)$ уже обладают свойством кучи
- Как переставить элементы поддерева с вершиной i , чтобы получилась куча?



Algorithm 4: Heapify(A,i)

```

1 l=Left(i);r=right(i);
2 if l < HeapSize ∧ A[l] > A[i] then
3   | largest=l
4 else
5   | largest=i
6 if r < HeapSize ∧ A[r] > A[largest] then
7   | largest=r
8 if largest ≠ i then
9   | swap(A[i],A[largest]);
10  | Heapify(A,largest)
  
```


Обеспечение основного свойства кучи ($A[\text{Parent}(i)] \succeq A[i], 0 < i < n$)

- Будем считать, что поддеревья с корнями $\text{Left}(i)$ и $\text{Right}(i)$ уже обладают свойством кучи
- Как переставить элементы поддерева с вершиной i , чтобы получилась куча?
- Худший случай: левое поддерево содержит вдвое больше элементов, чем правое
- В дереве из n элементов число узлов в левом и правом поддеревьях не превосходит $2n/3$

$$T(n) \leq T(2n/3) + \Theta(1) \Rightarrow T(n) = O(\log(n))$$

Algorithm 5: Heapify(A,i)

```

1 l=Left(i);r=right(i);
2 if l < HeapSize ∧ A[l] > A[i] then
3   | largest=l
4 else
5   | largest=i
6 if r < HeapSize ∧ A[r] > A[largest] then
7   | largest=r
8 if largest ≠ i then
9   | swap(A[i],A[largest]);
10  | Heapify(A,largest)

```

Построение кучи

Algorithm 6: BuildHeap(A)

```

1 HeapSize=size(A);
2 for ( $i = \lfloor (size(A) - 1)/2 \rfloor$ ;  $i \geq 0$ ;  $i--$ ) do
3    $\lfloor$  Heapify(A,i)

```

- Для $n = size(A) \leq 3$ корректность очевидна
- Вершины $n/2, \dots, n-1$ являются листьями, т.е. удовлетворяют свойству кучи
- Число вершин высоты h не превосходит $\lceil n/2^{h+1} \rceil$, высота кучи $\lfloor \log_2 n \rfloor$
- Время работы *Heapify* для вершины высоты h равно $O(h)$
- Время работы *BuildHeap* не превосходит

$$\sum_{h=0}^{\lfloor \log_2 n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O\left(n \sum_{h=0}^{\infty} \frac{h}{2^h}\right) =$$

$$O\left(n \frac{1/2}{(1 - 1/2)^2}\right) = O(n)$$

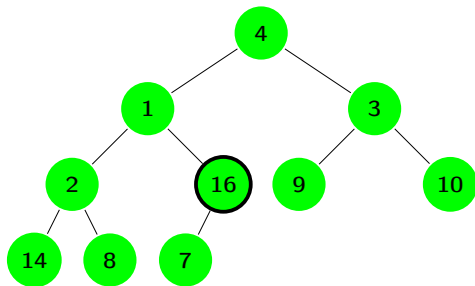
Построение кучи

Algorithm 7: BuildHeap(A)

```

1 HeapSize=size(A);
2 for ( $i = \lfloor (size(A) - 1)/2 \rfloor$ ;  $i \geq 0$ ;  $i--$ ) do
3    $\lfloor$  Heapify(A,i)
  
```

$A = (4, 1, 3, 2, 16, 9, 10, 14, 8, 7)$



- Для $n = size(A) \leq 3$ корректность очевидна
- Вершины $n/2, \dots, n-1$ являются листьями, т.е. удовлетворяют свойству кучи
- Число вершин высоты h не превосходит $\lceil n/2^{h+1} \rceil$, высота кучи $\lfloor \log_2 n \rfloor$
- Время работы *Heapify* для вершины высоты h равно $O(h)$
- Время работы *BuildHeap* не превосходит

$$\sum_{h=0}^{\lfloor \log_2 n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O\left(n \sum_{h=0}^{\infty} \frac{h}{2^h}\right) =$$

$$O\left(n \frac{1/2}{(1-1/2)^2}\right) = O(n)$$

Построение кучи

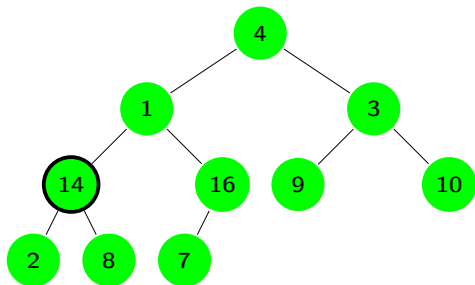
Algorithm 8: BuildHeap(A)

```

1 HeapSize=size(A);
2 for ( $i = \lfloor (size(A) - 1)/2 \rfloor$ ;  $i \geq 0$ ;  $i--$ ) do
3    $\lfloor$  Heapify(A,i)

```

$A = (4, 1, 3, 14, 16, 9, 10, 2, 8, 7)$



- Для $n = size(A) \leq 3$ корректность очевидна
- Вершины $n/2, \dots, n-1$ являются листьями, т.е. удовлетворяют свойству кучи
- Число вершин высоты h не превосходит $\lceil n/2^{h+1} \rceil$, высота кучи $\lfloor \log_2 n \rfloor$
- Время работы *Heapify* для вершины высоты h равно $O(h)$
- Время работы *BuildHeap* не превосходит

$$\sum_{h=0}^{\lfloor \log_2 n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O\left(n \sum_{h=0}^{\infty} \frac{h}{2^h}\right) =$$

$$O\left(n \frac{1/2}{(1 - 1/2)^2}\right) = O(n)$$

Построение кучи

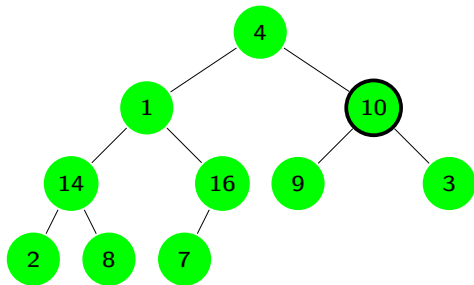
Algorithm 9: BuildHeap(A)

```

1 HeapSize=size(A);
2 for ( $i = \lfloor (size(A) - 1)/2 \rfloor$ ;  $i \geq 0$ ;  $i--$ ) do
3    $\lfloor$  Heapify(A,i)

```

$A = (4, 1, 10, 14, 16, 9, 3, 2, 8, 7)$



- Для $n = size(A) \leq 3$ корректность очевидна
- Вершины $n/2, \dots, n-1$ являются листьями, т.е. удовлетворяют свойству кучи
- Число вершин высоты h не превосходит $\lceil n/2^{h+1} \rceil$, высота кучи $\lfloor \log_2 n \rfloor$
- Время работы *Heapify* для вершины высоты h равно $O(h)$
- Время работы *BuildHeap* не превосходит

$$\sum_{h=0}^{\lfloor \log_2 n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O\left(n \sum_{h=0}^{\infty} \frac{h}{2^h}\right) =$$

$$O\left(n \frac{1/2}{(1 - 1/2)^2}\right) = O(n)$$

Построение кучи

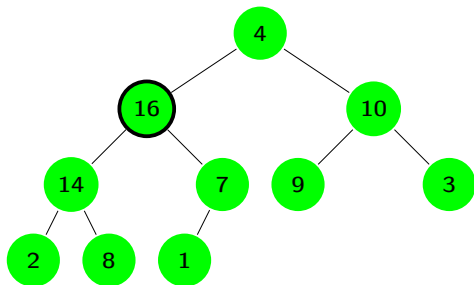
Algorithm 10: BuildHeap(A)

```

1 HeapSize=size(A);
2 for ( $i = \lfloor (size(A) - 1)/2 \rfloor ; i \geq 0 ; i --$ ) do
3    $\lfloor$  Heapify(A,i)

```

$A = (4, 16, 10, 14, 7, 9, 3, 2, 8, 1)$



- Для $n = size(A) \leq 3$ корректность очевидна
- Вершины $n/2, \dots, n-1$ являются листьями, т.е. удовлетворяют свойству кучи
- Число вершин высоты h не превосходит $\lceil n/2^{h+1} \rceil$, высота кучи $\lfloor \log_2 n \rfloor$
- Время работы *Heapify* для вершины высоты h равно $O(h)$
- Время работы *BuildHeap* не превосходит

$$\sum_{h=0}^{\lfloor \log_2 n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O\left(n \sum_{h=0}^{\infty} \frac{h}{2^h}\right) =$$

$$O\left(n \frac{1/2}{(1-1/2)^2}\right) = O(n)$$

Построение кучи

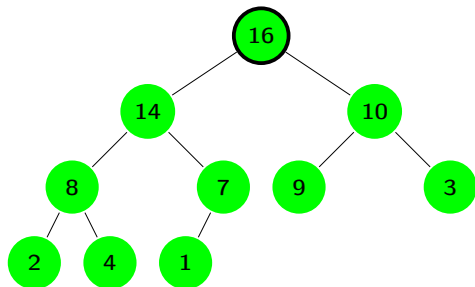
Algorithm 11: BuildHeap(A)

```

1 HeapSize=size(A);
2 for ( $i = \lfloor (size(A) - 1)/2 \rfloor$ ;  $i \geq 0$ ;  $i--$ ) do
3    $\lfloor$  Heapify(A,i)

```

$A = (16, 14, 10, 8, 7, 9, 3, 2, 4, 1)$



- Для $n = size(A) \leq 3$ корректность очевидна
- Вершины $n/2, \dots, n-1$ являются листьями, т.е. удовлетворяют свойству кучи
- Число вершин высоты h не превосходит $\lceil n/2^{h+1} \rceil$, высота кучи $\lfloor \log_2 n \rfloor$
- Время работы *Heapify* для вершины высоты h равно $O(h)$
- Время работы *BuildHeap* не превосходит

$$\sum_{h=0}^{\lfloor \log_2 n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O\left(n \sum_{h=0}^{\infty} \frac{h}{2^h}\right) =$$

$$O\left(n \frac{1/2}{(1 - 1/2)^2}\right) = O(n)$$

Производная геометрической прогрессии

$$\sum_{i=0}^t a^i = \frac{a^{t+1} - 1}{a - 1}$$

$$\sum_{i=0}^t i a^{i-1} = \frac{d}{da} \sum_{i=0}^t a^i = \frac{d}{da} \frac{a^{t+1} - 1}{a - 1} = \frac{(t+1)a^t(a-1) - (a^{t+1} - 1)}{(a-1)^2}$$

Сортировка с помощью кучи

- 1 Построим кучу
- 2 Поменяем местами корень дерева (максимальный элемент) и элемент в ячейке $n - 1$
- 3 Восстановим свойство кучи для начальных $n - 1$ элементов массива
- 4 Повторить п. 2–3 для первых $n - 1$ элементов массива

Сложность $O(n \log n)$

Algorithm 12: HeapSort(A)

```

1 BuildHeap(A)
2 for ( $i = n - 1; i > 0; i --$ ) do
3   swap(A[0], A[i])
4   HeapSize --
5   Heapify(A, 0)
```

Приоритетная очередь

Приоритетная очередь — структура данных, предназначенная для обслуживания множества $S = \{(key, value)\}$ с определенным на нем строгим слабым порядком, поддерживающая операции

- $Insert(S, x) : S \leftarrow S \cup \{x\}$
- $Maximum(S) : \text{возвращает } \max S$
- $ExtractMax(S) : \text{возвращает } \max S \text{ и удаляет соответствующий элемент из } S$
- $IncreaseKey(S, x, k) : \text{увеличивает значение ключа элемента } x \text{ путем его замены на } k, \text{ т.е. } x.key \leftarrow k, x.key \prec k$

Реализация приоритетной очереди на основе кучи

Algorithm 13: Maximum(*A*)

1 return *A*[0]

Сложность $\Theta(1)$

Algorithm 14: ExtractMax(*A*)

1 if *HeapSize* < 1 then

2 └ Error: очередь пуста

3 max=*A*[0];*A*[0]=*A*[*HeapSize*-1];

4 *HeapSize*=*HeapSize*-1;

5 Heapify(*A*,0);

6 return max

Сложность $O(\log(n))$

Algorithm 15: IncreaseKey(*A*,*i*,*k*)

1 if $k \prec A[i].key$ then

2 └ Error: Новый ключ меньше предыдущего

3 *A*[*i*].key=*k*;

4 while $i > 0 \wedge A[Parent(i)] \prec A[i]$ do

5 └ swap(*A*[*i*],*A*[*Parent*(*i*)]); *i*=*Parent*(*i*)

Сложность $O(\log n)$

Algorithm 16: Insert(*A*,(*key*,*value*))

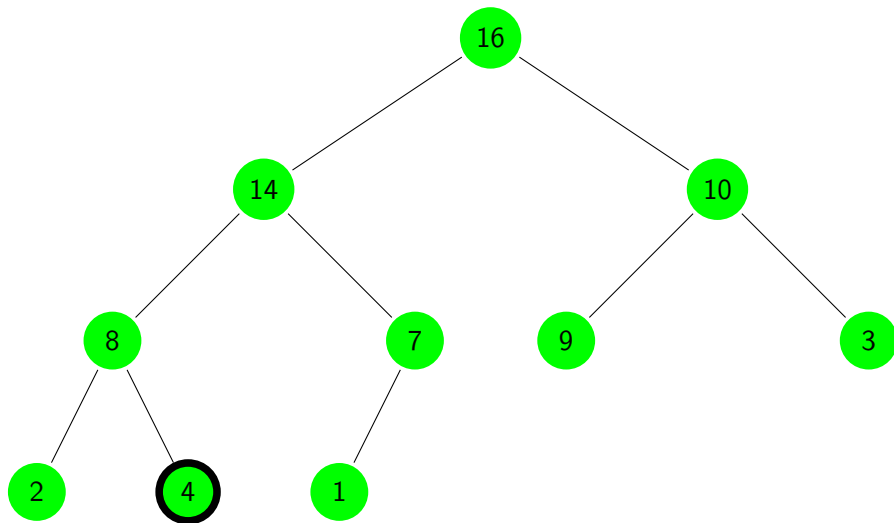
1 *A*[*HeapSize*]=($-\infty$, *value*);

2 *HeapSize*=*HeapSize*+1;

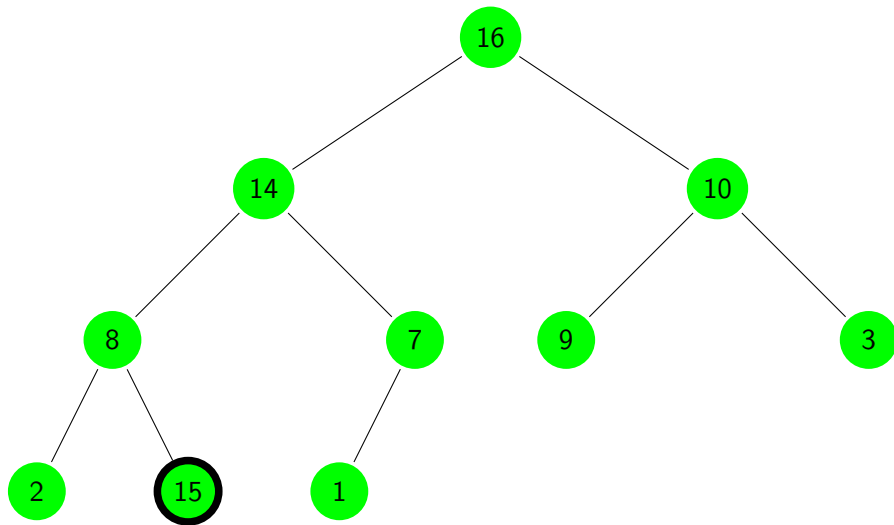
3 IncreaseKey(*A*,*HeapSize*-1,*key*)

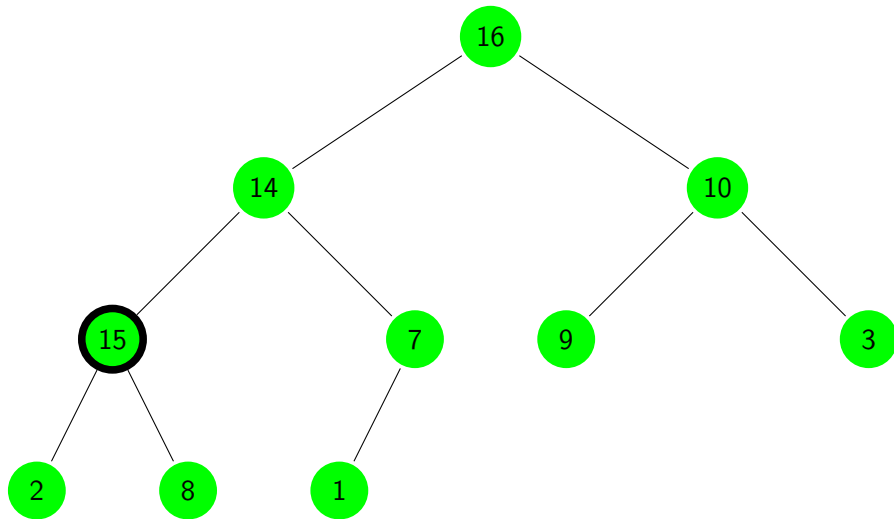
Сложность $O(\log n)$

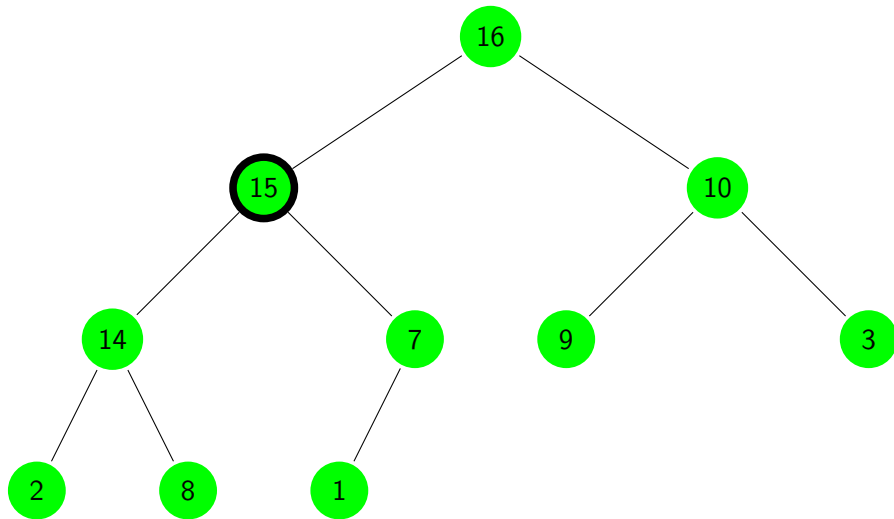
Пример работы *IncreaseKey*



Пример работы *IncreaseKey*



Пример работы *IncreaseKey*

Пример работы *IncreaseKey*

Быстрая сортировка

Сортировка массива $(A[p], \dots, A[r])$

- Переставим элементы массива так, чтобы $A[p], A[p+1], \dots, A[q] \preceq A[q+1], A[q+2], \dots, A[r]$, где $p \leq q < r$
- Вызовем рекурсивно процедуру сортировки для массивов $(A[p], \dots, A[q])$ и $(A[q+1], \dots, A[r])$

Как выбрать q ?

Algorithm 17: QuickSort(A, p, r)

```

1 if  $p < r$  then
2    $q = \text{Partition}(A, p, r)$ ;
3   QuickSort( $A, p, q$ );
4   QuickSort( $A, q+1, r$ );

```

Разбиение массива $A[p..r]$

	5	3	2	6	4	1	3	7	
$\uparrow i$									$\uparrow j$
	5	3	2	6	4	1	3	7	
$\uparrow i$							$\uparrow j$		
	5	3	2	6	4	1	3	7	
	$\uparrow i$						$\uparrow j$		
	3	3	2	6	4	1	5	7	
	$\uparrow i$						$\uparrow j$		
	3	3	2	6	4	1	5	7	
	$\uparrow i$					$\uparrow j$			
	3	3	2	6	4	1	5	7	
				$\uparrow i$		$\uparrow j$			
	3	3	2	1	4	6	5	7	
				$\uparrow i$		$\uparrow j$			
	3	3	2	1	4	6	5	7	
					$\uparrow j$	$\uparrow i$			

Algorithm 18: Partition(A, p, r)

```

1  $x = A[p]; i = p - 1; j = r + 1;$ 
2 while true do
3     repeat
4          $j = j - 1$ 
5     until  $A[j] \leq x;$ 
6     repeat
7          $i = i + 1$ 
8     until  $A[i] \geq x;$ 
9     if  $i < j$  then
10         $\text{swap}(A[i], A[j])$ 
11    else
12        return  $j$ 

```

Сложность $\Theta(r - p + 1)$

Корректность алгоритма разбиения ($p \leq r$)

- На строке 9 всегда выполняется $p \leq i, j \leq r$
 К первому проходу выполняется $p \leq j \leq r, i = p$
 В начале последующих итераций WHILE $p \leq i < j$
 На каждой итерации WHILE значение $A[p]$ не увеличивается
 $\Rightarrow p \leq j$ на строке 9
 После первой итерации $A[r] \succeq x \Rightarrow i \leq r$ на строке 9
- Если $p < r$, то на строке 12 $j < r$
- На строке 12 $A[s] \preceq A[t], p \leq s \leq j, j+1 \leq t \leq r$

Algorithm 19: Partition(A, p, r)

```

1   $x = A[p]; i = p - 1; j = r + 1;$ 
2  while true do
3      repeat
4           $j = j - 1$ 
5          until  $A[j] \preceq x;$ 
6          repeat
7               $i = i + 1$ 
8              until  $A[i] \succeq x;$ 
9          if  $i < j$  then
10              $\text{swap}(A[i], A[j])$ 
11         else
12             return  $j$ 
```

Время работы: крайние случаи

- Худший случай: на каждом шаге одна часть разбиения содержит 1 элемент, вторая — $n - 1$
 - $T(n) = \max_{1 \leq q \leq n-1} (T(q) + T(n - q)) + \Theta(n)$
 - Предположим, что
$$T(q) \leq cq^2 \Rightarrow T(n) \leq c \max_{1 \leq q < n} (q^2 + (n - q)^2) + \Theta(n) \leq c(1 + (n - 1)^2) + \Theta(n)$$
$$\Rightarrow T(n) \leq cn^2 - 2c(n - 1) + \Theta(n) \leq cn^2 \text{ (для достаточно большого } c \text{)}$$

Время работы: крайние случаи

- Худший случай: на каждом шаге одна часть разбиения содержит 1 элемент, вторая — $n - 1$
 - $T(n) = \max_{1 \leq q \leq n-1} (T(q) + T(n - q)) + \Theta(n)$
 - Предположим, что

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

$$\Rightarrow T(n) \leq cn^2 - 2c(n - 1) + \Theta(n) \leq cn^2 \text{ (для достаточно большого } c \text{)}$$
 - $T(n) = T(n - 1) + \Theta(n) = \sum_{k=1}^n \Theta(k) = \Theta(\sum_{k=1}^n k) = \Theta(n^2)$

Время работы: крайние случаи

- Худший случай: на каждом шаге одна часть разбиения содержит 1 элемент, вторая — $n - 1$
 - $T(n) = \max_{1 \leq q \leq n-1} (T(q) + T(n - q)) + \Theta(n)$
 - Предположим, что

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

$$\Rightarrow T(n) \leq cn^2 - 2c(n - 1) + \Theta(n) \leq cn^2 \text{ (для достаточно большого } c \text{)}$$
 - $T(n) = T(n - 1) + \Theta(n) = \sum_{k=1}^n \Theta(k) = \Theta(\sum_{k=1}^n k) = \Theta(n^2)$
 - Данный случай возникает, в частности, для отсортированных массивов (сортировка вставками при этом имеет сложность $\Theta(n)$)

Время работы: крайние случаи

- Худший случай: на каждом шаге одна часть разбиения содержит 1 элемент, вторая — $n - 1$
 - $T(n) = \max_{1 \leq q \leq n-1} (T(q) + T(n - q)) + \Theta(n)$
 - Предположим, что

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

$$\Rightarrow T(n) \leq cn^2 - 2c(n - 1) + \Theta(n) \leq cn^2 \text{ (для достаточно большого } c \text{)}$$
 - $T(n) = T(n - 1) + \Theta(n) = \sum_{k=1}^n \Theta(k) = \Theta(\sum_{k=1}^n k) = \Theta(n^2)$
 - Данный случай возникает, в частности, для отсортированных массивов (сортировка вставками при этом имеет сложность $\Theta(n)$)
- Лучший случай: на каждом шаге разбиение на две равные части

$$T(n) = 2T(n/2) + \Theta(n) = \Theta(n \log n)$$

Рандомизированный алгоритм

На практике вероятность появления отсортированных массивов существенно больше $2/n!$ \Rightarrow применяют рандомизированный алгоритм

Algorithm 20: RandPartition(A,p,r)

```

1 i=Random(p,r);
2 swap(A[p],A[i]);
3 return Partition(A,p,r)

```

Algorithm 21: RandQSort(A,p,r)

```

1 if  $p < r$  then
2   |   q=RandPartition(A,p,r);
3   |   RandQSort(A,p,q);
4   |   RandQSort(A,q+1,r);

```

Анализ среднего времени работы (рандомизированный алгоритм)

- Будем считать, что все элементы $(A[0], \dots, A[n-1])$ различны
- Значение, возвращаемое $Partition(A, 0, n-1)$, определяется числом элементов $A[i] \leq x$ ($\text{rank}(x)$), где $x = A[p]$. $P\{\text{rank}(x) = i\} = 1/n, 0 \leq i < n$
- $\text{rank}(x) = 0 \Rightarrow$ левая часть разбиения содержит 1 элемент (x)
- $\text{rank}(x) > 0 \Rightarrow$ левая часть содержит $\text{rank}(x)$ элементов
- Число элементов в левой части: $q = 1$ с вероятностью $P_1 = 1/n$, и $q : 1 < q < n$ с вероятностью $P_q = 1/n$
- Среднее время работы

$$\begin{aligned}
 T(n) &= \sum_{q=1}^{n-1} P_q (T(q) + T(n-q)) + \Theta(n) = \frac{1}{n} (T(1) + T(n-1)) + \frac{1}{n} \sum_{q=1}^{n-1} (T(q) + T(n-q)) + \Theta(n) = \\
 &= \frac{1}{n} O(n^2) + \frac{1}{n} \sum_{q=1}^{n-1} (T(q) + T(n-q)) + \Theta(n) = \frac{2}{n} \sum_{q=1}^{n-1} T(q) + \Theta(n) = O(n \log n)
 \end{aligned}$$

Решение рекуррентного уравнения

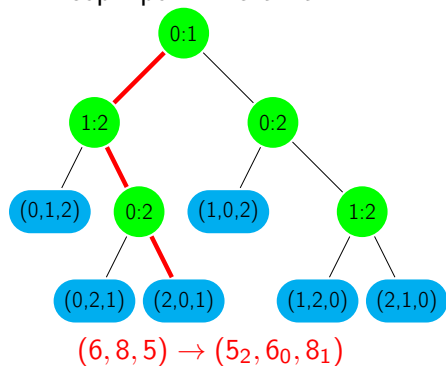
Предположим, что $T(n) \leq an \log n + b$. Это верно при $n = 1$ при достаточно больших a, b

$$\begin{aligned}
 T(n) &= \frac{2}{n} \sum_{q=1}^{n-1} T(q) + \Theta(n) \leq \frac{2}{n} \sum_{q=1}^{n-1} (aq \log q + b) + \Theta(n) = \frac{2a}{n} \sum_{q=1}^{n-1} q \log q + 2b \frac{n-1}{n} + \Theta(n) \leq \\
 &\leq \frac{2a}{n} \left(\sum_{q=1}^{\lceil \frac{n}{2} \rceil - 1} q \log q + \sum_{q=\lceil \frac{n}{2} \rceil}^{n-1} q \log q \right) + 2b + \Theta(n) \leq \\
 &\leq \frac{2a}{n} \left(\sum_{q=1}^{\lceil \frac{n}{2} \rceil - 1} q(\log n - 1) + \sum_{q=\lceil \frac{n}{2} \rceil}^{n-1} q \log n \right) + 2b + \Theta(n) \leq \\
 &\leq \frac{2a}{n} \left(\log n \sum_{q=1}^{n-1} q - \sum_{q=1}^{\lceil \frac{n}{2} \rceil - 1} q \right) + 2b + \Theta(n) \leq \frac{2a}{n} \left(\frac{n(n-1)}{2} \log n - \frac{1}{2} \left(\frac{n}{2} - 1 \right) \frac{n}{2} \right) + 2b + \Theta(n) \leq \\
 &\leq an \log n - an/4 + 2b + \Theta(n) = an \log n + b + (\Theta(n) + b - an/4) \leq an \log n + b
 \end{aligned}$$

Модель дерева решений

- Дерево решений — двоичное дерево, в котором представлены операции сравнения, выполняемые некоторым алгоритмом сортировки
- Внутренние узлы помечены $i : j, 0 \leq i, j < n$. Это означает сравнение $A[i], A[j]$
 - Левое поддерево: $A[i] \leq A[j]$
 - Правое поддерево: $A[i] > A[j]$
- Листья помечены перестановками $(\pi(0), \pi(1), \dots, \pi(n-1))$, задающими всевозможные упорядочения $A[\pi(0)] \leq A[\pi(1)] \leq \dots \leq A[\pi(n-1)]$
- Выполнение алгоритма сортировки — проход от корня дерева к одному из листьев
- Корректный алгоритм сортировки способен реализовать любую из $n!$ перестановок

Дерево решений для алгоритма сортировки вставками



Нижняя оценка сложности сортировки в худшем случае

Теорема

В наихудшем случае сложность любого алгоритма сортировки сравнением составляет $\Omega(n \log n)$ сравнений

Доказательство.

- Сортировка n элементов соответствует дереву решений высоты h с l достижимыми листьями
- Каждая из $n!$ перестановок соответствует одному из листьев $\Rightarrow n! \leq l$
- Число листьев в двоичном дереве высоты h удовлетворяет $l \leq 2^h$
- $n! \leq 2^h \Rightarrow h \geq \log(n!) = \log\left(\sqrt{2\pi n} \left(\frac{n}{e}\right)^n \left(1 + \Theta\left(\frac{1}{n}\right)\right)\right) = \Theta(n \log n)$



Следствие

Пирамидальная сортировка (сортировка на куче) и сортировка слиянием — асимптотически

Сортировка подсчетом

- Пусть ключи сортируемых элементов $A[i].key \in \mathcal{Z}, \mathcal{Z} = \{0, 1, 2, \dots, k\}, k = O(n)$
- Подсчитаем количество вхождений каждого числа из \mathcal{Z} в сортируемый массив
- Выпишем каждое число из \mathcal{Z} столько раз, сколько оно встретилось
- Сложность $\Theta(n + k)$; алгоритм устойчив

Algorithm 22: CountingSort(A,B,k)

```

1 for ( $i = 0; i \leq k; i++$ ) do
2    $C[i] = 0$ 
3 for ( $j = 0; j < \text{length}(A); j++$ ) do
4    $C[A[j].key]++$ 
5 for ( $i = 1; i \leq k; i++$ ) do
6    $C[i] += C[i-1]$ 
7 for ( $j = \text{length}(A) - 1; j \geq 0; j--$ ) do
8    $C[A[j]]--$ ;
9    $B[C[A[j].key]] = A[j];$ 

```

$A = (2, 5, 3, 0, 2, 3, 0, 3)$

$C = (2, 2, 4, 7, 7, 8)$

$B = (0, 0, 2, 2, 3, 3, 3, 5)$

Сортировка подсчетом

- Пусть ключи сортируемых элементов $A[i].key \in \mathcal{Z}, \mathcal{Z} = \{0, 1, 2, \dots, k\}, k = O(n)$
- Подсчитаем количество вхождений каждого числа из \mathcal{Z} в сортируемый массив
- Выпишем каждое число из \mathcal{Z} столько раз, сколько оно встретилось
- Сложность $\Theta(n + k)$; алгоритм устойчив
- Это не противоречит вышедшей теореме
 - Алгоритм не использует сравнения

Algorithm 23: CountingSort(A, B, k)

```

1 for ( $i = 0; i \leq k; i++$ ) do
2    $C[i] = 0$ 
3 for ( $j = 0; j < \text{length}(A); j++$ ) do
4    $C[A[j].key]++$ 
5 for ( $i = 1; i \leq k; i++$ ) do
6    $C[i] += C[i-1]$ 
7 for ( $j = \text{length}(A) - 1; j \geq 0; j--$ ) do
8    $C[A[j]]--$ ;
9    $B[C[A[j].key]] = A[j];$ 

```

 $A = (2, 5, 3, 0, 2, 3, 0, 3)$
 $C = (2, 2, 4, 7, 7, 8)$
 $B = (0, 0, 2, 2, 3, 3, 3, 5)$

Сортировка подсчетом

- Пусть ключи сортируемых элементов $A[i].key \in \mathcal{Z}, \mathcal{Z} = \{0, 1, 2, \dots, k\}, k = O(n)$
- Подсчитаем количество вхождений каждого числа из \mathcal{Z} в сортируемый массив
- Выпишем каждое число из \mathcal{Z} столько раз, сколько оно встретилось
- Сложность $\Theta(n + k)$; алгоритм устойчив
- Это не противоречит вышедшей теореме
 - Алгоритм не использует сравнения
 - Если $n > 2^h$, где h — разрядность ЭВМ, то для реализации операций с $C[i]$ понадобится длинная арифметика с числом разрядов $\eta = \lceil \log_h(n) \rceil$ и сложностью операций $\Theta(\eta) \Rightarrow$ сложность сортировки $\Theta(n \log n)$

Algorithm 24: CountingSort(A,B,k)

```

1 for ( $i = 0; i \leq k; i++$ ) do
2    $C[i] = 0$ 
3 for ( $j = 0; j < \text{length}(A); j++$ ) do
4    $C[A[j].key]++$ 
5 for ( $i = 1; i \leq k; i++$ ) do
6    $C[i] += C[i-1]$ 
7 for ( $j = \text{length}(A) - 1; j \geq 0; j--$ ) do
8    $C[A[j]]--$ ;
9    $B[C[A[j].key]] = A[j]$ ;
```

$A = (2, 5, 3, 0, 2, 3, 0, 3)$

$C = (2, 2, 4, 7, 7, 8)$

$B = (0, 0, 2, 2, 3, 3, 3, 5)$

Поразрядная сортировка

- Пусть даны n натуральных чисел $\sum_{j=0}^{d-1} a_{ij}b^j, 0 \leq a_{ij} < b$
- Отсортируем (устойчиво!) их по младшей цифре a_{i0} , объединим в массив
- Отсортируем (устойчиво!) их по a_{i1} , объединим в массив
- ...
- Корректность: доказательство индукцией по d
- Если сортировка по каждой цифре имеет сложность $\Theta(n + b)$, то алгоритм имеет сложность $\Theta(d(n + b))$

3 2 9	7 2 0	7 2 0	3 2 9
4 5 7	3 5 5	3 2 9	3 5 5
6 5 7	4 3 6	4 3 6	4 3 6
8 3 9	→ 4 5 7	→ 8 3 9	→ 4 5 7
4 3 6	6 5 7	3 5 5	6 5 7
7 2 0	3 2 9	4 5 7	7 2 0
3 5 5	8 3 9	6 5 7	8 3 9

Лемма

Пусть даны n чисел разрядности h и число $r \leq h$. Их поразрядная сортировка имеет сложность $\Theta((b/r)(n + 2^r))$

Карманная сортировка (сортировка вычерпыванием)

- Рассмотрим сортировку значений x , равномерно распределенных на $[0, 1)$
- Разложим сортируемые числа по карманам (ведрам, ящикам, ...), соответствующим интервалам $[\frac{i}{n}, \frac{i+1}{n})$, $0 \leq i < n$
- Карман — список
- Отсортируем числа в каждом кармане
- Последовательно выгрузим содержимое всех карманов

Algorithm 25: BucketSort(A)

```

1 n=length(A);j=0;
2 for (i = 0; i < n; i++) do
3   Append( $B_{\lfloor nA[i] \rfloor}$ , A[i]);
4 for (i = 0; i < n; i++) do
5   Сортировка вставками  $B_i$ ;
6   Переместить содержимое  $B_i$  в
    $A[j], \dots, A[j + |B_i| - 1]$ ;
7    $j := j + |B_i|$ 

```

$A = (0.78, 0.17, 0.39, 0.26, 0.72, 0.94, 0.21, 0.12, 0.23, 0.68)$

i	0	1	2	3	4	5	6	7	8	9
B_i		0.12	0.21	0.39			0.68	0.72		0.94
		0.17	0.23				0.78			

Сложность карманной сортировки

- Сложность сортировки кармана $O(n_i^2)$, где n_i — число элементов в B_i
- Матожидание суммарной сложности сортировки

$$M[C] = \sum_{i=0}^{n-1} O(M[n_i^2]) = O\left(\sum_{i=0}^{n-1} M[n_i^2]\right)$$

- Вероятность попадания $A[j]$ в $B[i]$ равна $p = 1/n$
- $P\{n_i = k\} = C_n^k p^k (1-p)^{n-k}, 0 \leq k \leq n$
- $M[n_i] = np = 1$. $D[n_i] = np(1-p) = 1 - 1/n$

$$M[n_i^2] = D[n_i] + (M[n_i])^2 = 1 - 1/n + 1 = 2 - 1/n = \Theta(1)$$

- Средняя сложность сортировки $M[C] = O(n)$

Минимум и максимум

Algorithm 26: Minimum(A)

```

1 min=A[0];
2 for (i=1;i<length(A);i++) do
3   if A[i] < min then
4     min=A[i]
5 return min

```

- Нахождение минимума (максимума) требует ровно $n - 1$ сравнений
- $n - 1$ элементов должны “проиграть” минимальному

Algorithm 27: MinimumAndMaximum(A)

```

1 min = A[0]; max = A[1 - (length(A) mod 2)];
2 if max < min then
3   swap(min,max)
4 for (i = 2 - (length(A) mod 2); i < length(A); i += 2) do
5   a=A[i];b=A[i+1];
6   if b < a then
7     swap(a,b)
8   if a < min then
9     min=a
10  if max < b then
11    max=b
12 return (min,max)

```

$3 \lceil n/2 \rceil - 2$ сравнений

Порядковые статистики

- i -ая порядковая статистика множества (массива) A — его i -ый элемент в порядке возрастания, $0 \leq i < |A|$
 - $\text{ord}_2(30, 3, 5, 10) = 10$
- Минимум — $\text{ord}_0(A)$
- Максимум — $\text{ord}_{|A|-1}(A)$
- Медиана:
$$\begin{cases} \text{ord}_{(|A|-1)/2}(A), & |A| \text{ нечетная} \\ \{ \underbrace{\text{ord}_{|A|/2-1}(A)}_{\text{нижняя медиана}}, \underbrace{\text{ord}_{|A|/2}(A)}_{\text{верхняя медиана}} \}, & |A| \text{ четная} \end{cases}$$
- Далее будем рассматривать нижнюю медиану
- Выбор i -ой порядковой статистики может быть выполнен со сложностью $O(n \log n)$ посредством сортировки, но это неэффективно

Выбор i -ой порядковой статистики с линейной сложностью

Algorithm 28: Partition(A,p,r)

```

1 x=A[p];i=p-1;j=r+1;
2 while true do
3   repeat
4     | j=j-1
5   until A[j] ≤ x;
6   repeat
7     | i=i+1
8   until A[i] ≥ x;
9   if i < j then
10    | swap(A[i],A[j])
11  else
12    | return j

```

Сложность $\Theta(r - p + 1)$

Algorithm 29: RandPartition(A,p,r)

```

1 i=Random(p,r);
2 swap(A[p],A[i]);
3 return Partition(A,p,r)

```

Algorithm 30: RandSelect(A,p,r,i)

```

1 if p=r then
2   | return A[p]
3 q=RandPartition(A,p,r);
4 k=q-p+1;
5 if i=k then
6   | return A[q]
7 else if i < k then
8   | return RandSelect(A,p,q-1,i)
9 else
10  | return RandSelect(A,q+1,r,i-k)

```

Выбор i -ой порядковой статистики с линейной сложностью

- *RandPartition* разбивает массив на подмассивы $A[p..q-1] \prec A[q+1..r]$
- $A[q]$ — опорный элемент
- Строка 5: не является ли $A[q]$ искомым значением?
- Строка 8: искомое значение в подмассиве $A[p..q-1]$
- Строка 10: искомое значение в подмассиве $A[q+1..r]$

Algorithm 31: RandSelect(A, p, r, i)

```

1 if  $p=r$  then
2   | return  $A[p]$ 
3  $q = \text{RandPartition}(A, p, r)$ ;
4  $k = q - p + 1$ ;
5 if  $i = k$  then
6   | return  $A[q]$ 
7 else if  $i < k$  then
8   | return  $\text{RandSelect}(A, p, q-1, i)$ 
9 else
10  | return  $\text{RandSelect}(A, q+1, r, i-k)$ 

```

Выбор i -ой порядковой статистики с линейной сложностью: сложность

- Худший случай: *RandPartition* всегда возвращает номер наибольшего элемента
 $\Rightarrow \Theta(n^2), n = r - p + 1$
- $P\{q - p + 1 = k\} = 1/n, 1 \leq k \leq n$
- $\chi_k = \begin{cases} 1, & q - p + 1 = k \\ 0, & \text{иначе} \end{cases}$
- $M[\chi_k] = 1/n$
- Предположим, что искомый элемент всегда попадает в самую большую часть разбиения

$$T(n) \leq \sum_{k=1}^n \chi_k (T(\max(k-1, n-k)) + O(n))$$

Algorithm 32: RandSelect(A,p,r,i)

```

1  if p=r then
2    return A[p]
3  q=RandPartition(A,p,r);
4  k=q-p+1;
5  if i=k then
6    return A[q]
7  else if i < k then
8    return RandSelect(A,p,q-1,i)
9  else
10   return RandSelect(A,q+1,r,i-k)

```

Выбор i -ой порядковой статистики с линейной сложностью: сложность

$$T(n) \leq \sum_{k=1}^n \chi_k T(\max(k-1, n-k)) + O(n)$$

$$\begin{aligned} M[T(n)] &\leq \sum_{k=1}^n M[\chi_k T(\max(k-1, n-k))] + O(n) = \sum_{k=1}^n M[\chi_k] M[T(\max(k-1, n-k))] + O(n) = \\ &= \sum_{k=1}^n \frac{1}{n} M[T(\max(k-1, n-k))] + O(n) \end{aligned}$$

- χ_k и $T(\max(k-1, n-k))$ — независимые случайные величины
- $\max(k-1, n-k) = \begin{cases} k-1, & k > \lceil n/2 \rceil \\ n-k, & k \leq \lceil n/2 \rceil \end{cases} \Rightarrow$ для четного n каждое слагаемое $T(\lceil n/2 \rceil), \dots, T(n-1)$ появляется дважды

Выбор i -ой порядковой статистики с линейной сложностью: сложность

$$M[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} M[T(k)] + O(n)$$

Предположим, что $M[T(n)] \leq cn$, а для малых n : $T(n) = O(1)$

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

Нужно показать, что при больших n : $cn/4 - c/2 - an \geq 0 \Leftrightarrow n(c/4 - a) \geq c/2$. $c > 4a \Rightarrow n \geq \frac{c/2}{c/4 - a}$

Заключение

- Сортировка может быть выполнена со сложностью $O(n \log n)$
- Иногда сортировка может быть выполнена со сложностью $O(n)$
- Порядковые статистики могут быть найдены со сложностью $O(n)$