

Algorithmic Design

LEZIONI	LUN	MAR	MER	GIO
9 - 11	alg design		stat ML	alg design
11 - 13	stat methods	stat methods	stat ML	stat ML
14 - 16	big data	big data		stat ML
16 - 18	big data	big data		stat methods

Alberto Casagrande
acasagrande@units.it

ricevimento: room 330 h2bis, monday and tuesday, 11-12

book: *Introduction to algorithms*, Cormen

[moodle](#)

key: Dijkstra

05/03/18

Algorithms

low level descriptions of well defined computational procedures meant to be automatized.

Computational model = single processor single thread model (RAM)

- infinite memory
- infinite variables
- basic algebraic operations
- basic algebraic relations

We will assume constant execution time for each of the above operations.

pseudo-code rules:

- if then, while, for, repeat statements
- \leftarrow denotes assignments
- blocks are identified by indentation
- array indexes begin from 1 (not 0)
- parameters are passed by value

Big-o notation

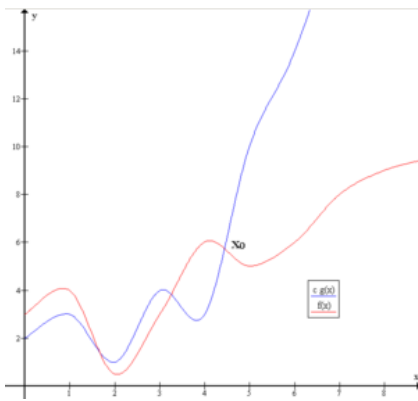
speedup theorem

given any real $c > 0$ and any Turing machine solving a problem in time $f(n)$, there is another machine that solves the same problem in time at most $cf(n) + n + 2$.

This means that constants do not matter, so we end up considering asinhtotic time complexity.

$$\text{def. } O(f) = \{g | \exists c > 0, n_0 \in \mathbb{N} : \forall n \geq n_0 |g(n)| \leq c|f(n)|\}$$

for example, $n + 1 \in O(n) \subset O(n^2) \subset O(n^{n!})$



def. $\Omega(f) = \{g | \exists c > 0, n_0 \in \mathbb{N} : \forall n \geq n_0 |g(n)| \geq c|f(n)|\}$
 $\Theta(f) = O(f) \cap \Omega(f)$

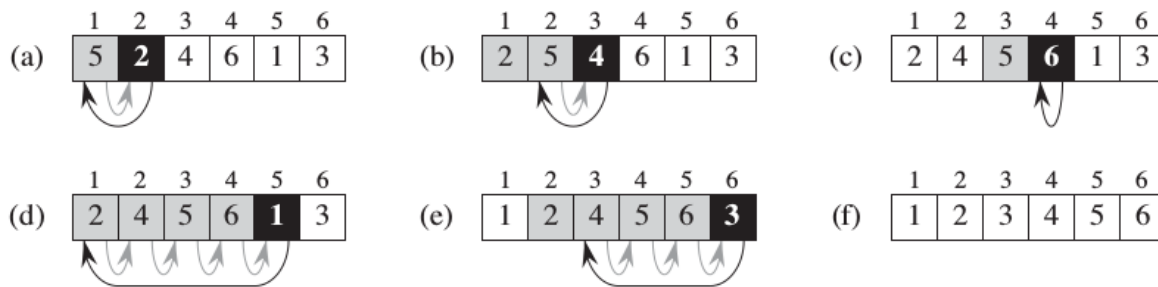
Sorting algorithms

p.26 Cormen

input is an array A

output is A such that for each $i \in [1, |A|)$, $A[i] \leq A[i + 1]$

Insertion sort



insertion sort

Can take different amounts of time to sort two input sequences of the same size depending on how nearly sorted they already are.

```

1 insertion_sort(A)
2   for j ← 2 to |A| do           # |A|-1 times
3     k ← A[j]                     # O(1)
4     i ← j - 1                   # O(1)
5     while i > 0 and A[i] > k do  # at most j-1
6       A[i + 1] ← A[i]           # O(1)
7       i ← i - 1                 # O(1)
8     A[i + 1] ← k                # O(1)

```

Let's compute the complexity of this algorithm:

$$\begin{aligned}
 T(|A|) &= \sum_{j=2}^{|A|} \sum_{i=1}^{j-1} (O(1) + O(1)) = \sum_{j=2}^{|A|} \sum_{i=1}^{j-1} O(1) \\
 &= \sum_{j=2}^{|A|} O(j-2) = \sum_{j=2}^{|A|} O(j) = O\left(\frac{|A|(|A|+1)}{2}\right) = O(|A|^2)
 \end{aligned}$$

Class	Sorting algorithm
Data structure	Array
Worst-case performance	$O(n^2)$ comparisons, swaps
Best-case performance	$O(n)$ comparisons, $O(1)$ swaps
Average performance	$O(n^2)$ comparisons, swaps
Worst-case space complexity	$O(n)$ total, $O(1)$ auxiliary

insertion sort summary

Select problem

→ p.220, 302-308 Cormen

→ [time_bounds_selection.pdf](#)

Given an array A we want to find the j -th element of the sorted A , without sorting it.

- input: an unsorted array A and an index $j \in [1, |A|]$
- output: a value $v \in A$ such that A contains $j - 1$ elements smaller than v

example

input: $A = [5, 2, 4, 6, 1, 3]$ and $j = 5$

output: 5

Quickselect algorithm

Median of the medians

Since the complexity of sorting is $O(n \log n)$, we aim at having a faster algorithm, maybe $O(n)$.

A possible strategy is the following one:

- split A in $\lceil \frac{n}{5} \rceil$ blocks of dimension 5 in $O(n)$
- compute the block medians in $O(n)$
- compute the median of all the medians (pivot) recursively
- compute $[k_1, k_2]$ such that for all $k \in [k_1, k_2]$ $A[k]$ is equal to the median of the medians
- if $j \in [k_1, k_2]$ then $A[j]$ is the median of the medians, otherwise:
 - if $j < k_1$ go back to 4. using $[0, k_1)$
 - if $j > k_2$ go back to 4. using $(k_2, |A|]$

alternative definition:

- given the interval $[k_1, k_2]$, compute $k \in [k_1, k_2]$ such that $A[k]$ is equal to the median of the medians
- if $j = k$ then $A[j]$ is the median of the medians, otherwise:
 - if $j < k$ go back to 4. using $[0, k - 1]$
 - if $j > k$ go back to 4. using $[k + 1, |A|]$

One iteration on a randomized set of 100 elements from 0 to 99

	12	15	11	2	9	5	0	7	3	21	44	40	1	18	20	32	19	35	37	39
	13	16	14	8	10	26	6	33	4	27	49	46	52	25	51	34	43	56	72	79
Medians	17	23	24	28	29	30	31	36	42	47	50	55	58	60	63	65	66	67	81	83
	22	45	38	53	61	41	62	82	54	48	59	57	71	78	64	80	70	76	85	87
	96	95	94	86	89	69	68	97	73	92	74	88	99	84	75	90	77	93	98	91

median of medians

The selection algorithm becomes:

```
1 select(A, j, begin=1, end=|A|)
2     // reorder the array
3     // sto supponendo che j sia 140...
4     if (end-begin+1 < 140)                // O(1)
5         insertion_sort(A, begin, end)      // O(|A|)
6         return A[j]
7     // find the median of medians
8     pivot = select_pivot(A,begin,end)      // O(|A|/5) + T(|A|/5)
9     // tri_partition partitions A[] in three parts:
10    // 1) A[begin,...,k_1-1] contains all elements smaller than the pivot
11    // 2) A[k_1,..., k_2] contains all occurrences of the pivot
12    // 3) A[k_2+1,...,end] contains all elements greater than the pivot
13    // finally it returns k_1 and k_2
14    (k_1, k_2) = tri_partition(A,begin,end,pivot)    // O(|A|)
```

```

15 // recursively call select in the correct interval
16 if  $j < k_1$ 
17     return select(A, j, begin,  $k_1 - 1$ )
18 if  $j > k_2$ 
19     return select(A, j,  $k_2 + 1$ , end)
20 return A[j]

```

some observations on the algorithm:

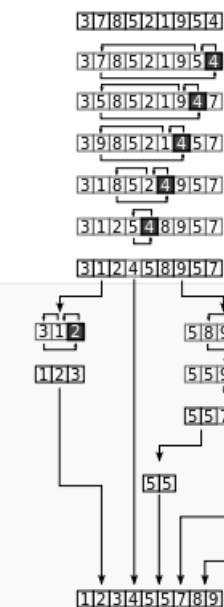
- `select_pivot` finds the pivot, which is the median of the block of medians.
- `tri_partition` takes all the elements in A smaller than the pivot at the begin of A and those greater at end
- Median of medians can also be used as a pivot strategy in [quicksort](#), yielding an optimal algorithm, with worst-case complexity $O(n \log n)$.
- in this case the complexity of insertion sort is $O(1)$ because we used the const value 140.

Select pivot algorithm

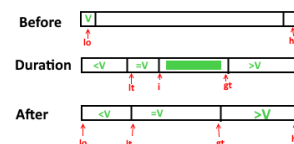
```

select_pivot(A, begin, end)
    num_of_blocks <- (end-begin+1)/5      # O(1)
    # initialize medians array
    medians <- array[num_of_blocks]      # O(|A|)
    # compute the medians for each block
    for i in [0, num_of_blocks]          # O(|A|/5) = O(|A|)
        # beginning point for each block
        cbegin = begin + 5 * i           # O(1)
        # reorder the block
        insertion_sort(A, cbegin, cbegin+4) # O(1)
        # add the i-th median to the array
        medians[i+1] <- A[cbegin+2]      # O(1)
    # find the median of medians
    return select(medians, j, num_of_blocks, (1+num_of_blocks)/2) # T(|A|/5)

```



quicksort



3 - way partitioning overview

the whole complexity of `select_pivot` is

$$T(n) = T(\lceil n/5 \rceil) + O(n/5) = T(\lceil n/5 \rceil) + O(n).$$

Let x be the pivot (median of medians) and notice that:

- at least half of the medians for each block are $\geq x$
- the blocks are $\lceil \frac{n}{5} \rceil$
- for each block there are at least 3 elements $\geq x$
- you have to exclude only two cases: the block containing x and the block having less than 5 elements (if 5 does not exactly divide n)

Considering all of this, the number of elements greater than pivot is at least

$$3(\lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil - 2) \geq \frac{3}{10}n - 6. \text{ The same inequality holds for the number of elements lower than the pivot.}$$

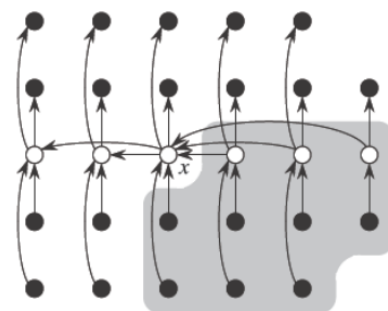
So the complexity for `select` is

$$T(n) = \begin{cases} O(1) & n < 140 \\ T(\lceil \frac{n}{5} \rceil) + T(\frac{7}{10}n + 6) + O(n) & n \geq 140 \end{cases}$$

Theorem. $T(n) \in O(n)$

proof. (by induction) p.222 Cormen

We want to prove that there exists a $c > 0$ such that $T(n) \leq cn \quad \forall n > 0$.



Let's choose c such that $T(n) \leq cn \forall n < 140$ and a s.t. $O(n) \leq an \forall n > 0$. For each $n \geq 140$ we have

$$\begin{aligned} T(n) &\leq c \lceil \frac{n}{5} \rceil + c(\frac{7}{10}n + 6) + an \\ &\leq cn/5 + c + 7cn/10 + 6c + an \\ &\leq cn + (-cn/10 + 7c + an) \end{aligned}$$

If the second portion is ≤ 0 we get $T(n) \leq cn$.

This is true because $n \geq 140 > 70$, so the inequality holds iff $c \geq 10an/(n - 70)$. If we choose $c \geq 20a$, noticing that $n/(n - 70) \geq 2$ we get the inequality. ■

If the complexity of tri_partition is $O(n)$ then the complexity of elect is $O(n)$.

12/03/18

Strassen's algorithm

→ p.735-741 Cormen

Given two matrices A and B, we can divide both in 4 parts, which in total are 8 blocks:

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \quad B = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$$

the method:

1. Split A and into 8 submatrices;
2. Compute 10 intermediate matrices by summing or subtracting submatrices from step 1. (S_1, S_2, \dots, S_{10});
3. Compute 7 matrices: P_1, \dots, P_7 , by multiplying S_1, \dots, S_{10} ;
4. Compute C_{11}, \dots, C_{22} by summing/subtracting P_1, \dots, P_7 .

The complexity cost is:

1. $O(1)$ using indexes
2. for each intermediate matrix is: $O(\frac{n^2}{2})$, which is equivalent to $O(n^2)$;
3. 7 times the cost of multiplying each sub-matrix: $7T(\frac{n}{2})$;
4. $O(n^2)$

where

$$T(n) = \begin{cases} O(1) & n > 1 \\ 7T(\frac{n}{2}) + O(n^2) & n \leq 1 \end{cases}$$

Setting $a = 7, b = 2$ and using master theorem we get $T(n) \in O(n^{\log_2 7})$.

The matrices are:

$$S_1 = B_{12} - B_{22}$$

$$S_2 = A_{11} + A_{12}$$

$$S_3 = A_{21} + A_{22}$$

$$S_4 = B_{21} - B_{11}$$

$$S_5 = A_{11} + A_{22}$$

$$S_6 = B_{11} + B_{22}$$

$$S_7 = A_{12} - A_{22}$$

$$S_8 = B_{21} + B_{22}$$

$$S_9 = A_{11} - A_{21}$$

$$S_{10} = B_{11} + B_{12}$$

$$P_1 = A_{11}S_1$$

$$P_2 = S_2B_{22}$$

$$P_3 = S_3B_{11}$$

$$P_4 = A_{22}S_4$$

$$P_5 = S_5S_6$$

$$P_6 = S_7S_8$$

$$P_7 = S_9S_{10}$$

$$C_{11} = P_5 + P_4 - P_2 + P_6$$

$$C_{12} = P_1 + P_2$$

$$C_{21} = P_3 + P_4$$

$$C_{22} = P_5 + P_1 - P_3 - P_7$$

☐ prove the equations and implement the code with matrices of size $\neq 2^n$

15/03/18

Graph algorithms

Strongly connected components

We have two main representations for graphs (V, E) :

- **Adjacency list:** map from nodes to the list of their ~~successors~~ adjacent nodes

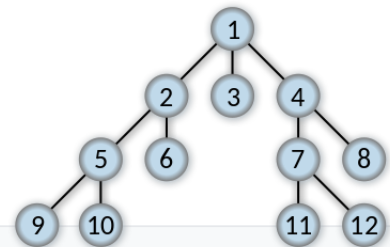
$Adj[v] \rightarrow$ list of nodes $w \in V$ such that $(v, w) \in E$

- **Adjacency matrix:** 0/1 matrix whose rows/columns correspond to the nodes of the graph. In position (i, j) we have 1 iff $(i, j) \in E$.

Breadth first search

\rightarrow p.531-539 Cormen

- Used for distance
- it uses a queue (FIFO)



BFS

```

1  BFS(G,v)                                # O(|V|+|E|)
2    for w in V do G = (V,E)               # O(|V|)
3      color[w] <- white
4      d[w] <- ∞
5      π[w] <- NIL
6  d[v] <- 0                                # O(1)
7  color[v] <- grey                         # O(1)
8  Enqueue(Q,v)                             # O(1)
9  while |Q| ≠ 0 do                         # O(|E|)
10     for w in Adj[head(Q)]
11       if color(w) == white then
12         d[w] <- d[head(Q)]+1
13         π[w] <- head(Q)
14         color[w] <- grey
15         Enqueue(Q,w)
16  color[head(Q)] <- black
17  Dequeue(Q)

```

So total time complexity is $O(|V| + |E|)$, since every vertex and every edge will be explored in the worst case. Notice that $O(|E|)$ may vary between $O(1)$ and $O(|V|^2)$ depending on how sparse the input graph is.

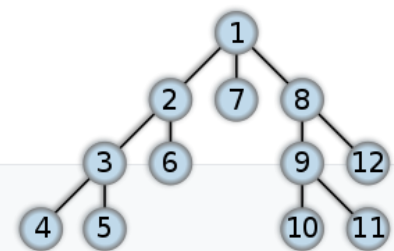
Depth first search

\rightarrow p.540-549 Cormen

Used for connectivity

We will have:

- $d : V \rightarrow \mathbb{N}$ discovery tree (preordering)
- $f : V \rightarrow \mathbb{N}$ final processing tree (postordering)
- $\pi : V \rightarrow \mathbb{N}$



DFS

```

1 DFS(G)                                     #  $O(|V| + |E|)$ 
2   for w in V do                             #  $O(|V|)$ 
3     color[w] <- white
4      $\pi[w]$  <- nil
5     d[w] <-  $\infty$ 
6     f[w] <-  $\infty$ 
7   time <- 0
8   for v in V do                             #  $O(|V|)$ 
9     if color[v] == white
10      time <- DFS_real(G,v,time)

```

where

```

1 DFS_real(G,v,time)
2   color[v] <- grey
3   d[v] <- time
4   time = time + 1
5   for w in Adj[v] do                       #  $O(|Adj[v]|)$ 
6     if color[w] == white
7        $\pi[w]$  <- v
8       time <- DFS_real(G,w,time)
9   f[v] <- time
10  time <- time + 1
11  return time

```

and the total complexity is $O(|V| + |E|)$.

Theorem. Let us consider the two intervals $I_v = [d[v], f[v]]$ and $I_w = [d[w], f[w]]$. Either:

- $I_v \cap I_w = \emptyset$
- $I_v \subseteq I_w$ and w is an ancestor of v in the DFS tree
- $I_w \subseteq I_v$ and v is an ancestor of w in the DFS tree

Tarjan's algorithm for SCC

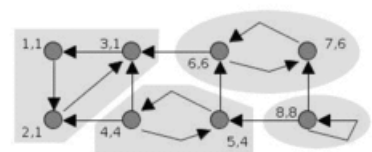
→ [p.615 Cormen](#)

→ [tarjan.pdf](#)

It is a sort of DFS + a map lowlink that remembers the smallest $d[v]$ among all the w reachable from v .

d = discovery time

f = final processing time



Tarjan's Algorithm Animation.gif

The algorithm takes a **directed graph** as input, and produces a **partition** of the graph's **vertices** into the graph's strongly connected components. Each vertex of the graph appears in exactly one of the strongly connected components. Any vertex that is not on a directed cycle forms a strongly connected component all by itself.

Each node v is assigned a unique integer $d[v]$, which numbers the nodes consecutively in the order in which they are discovered. It also maintains a value $lowlink[v]$ that represents the smallest index of any node known to be reachable from v through v 's DFS subtree, including itself. Therefore if $lowlink[v] = d[v]$ at the end of the recursive call over v , all the nodes reachable from v (e.g. those in the same SCC) have already been completed.

```

1 Tarjan_SCC(G)
2   # initialization
3   for v in V      # G(V,E)
4       d[v] <- ∞   # f[v] <- ∞
5       lowlink[v] <- ∞
6       color[v] <- white # v has not been visited
7   Q <- [] # successors queue
8   SCC <- []
9   time <- 0
10  for v in V
11      if color[v] == white
12          time <- Tarjan_SCC_real(G,v,time,Q,SCCs)
13  return SCCs

```

Next time we'll define `Tarjan_SCC_real`.

19/03/18

```

1 Tarjan_SCC_real(G,v,time,Q,SCCs)
2   # update v infos
3   d[v] <- time
4   lowlink[v] <- time
5   time <- time + 1
6   color[v] <- grey
7   Enqueue(Q,v)
8   for w in Adj[v] do
9       if color[w] == white # w has not been visited
10          time <- Tarjan_SCC_real(G,w,time,Q,SCCs)
11          lowlink[v] <- min(lowlink[v], lowlink[w])
12       else
13          if color[w] == grey # w has already been visited
14              lowlink[v] <- min(lowlink[v], lowlink[w])
15   if lowlink[v] == d[v] then
16       # v is a representative for a new SCC
17       SCC <- []
18       while Head(Q) != v do
19           SCC.append(Head(Q))
20           Dequeue(Q)
21       SCC.append(SCC)
22   return time

```

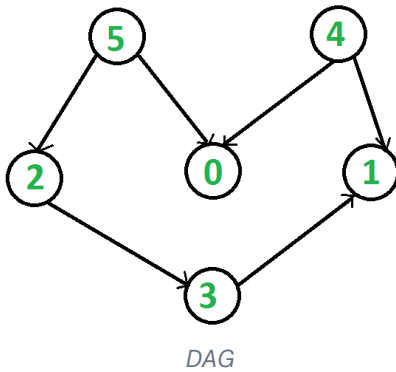
The whole complexity of `tarjan_SCC_real` is $O(|V|)$, so the complexity of `tarjan_SCC` is $O(|V| + |E|)$.

We cannot perform better than this algorithm for SCCs, because it has the same complexity of reading the whole graph.

Topological Sort

→ p. 612 Cormen

Definition of Strongly Connected Components. Tarjan's algorithm to compute the strongly connected components of a graph. Intuition, code, and complexity analysis. Definition of topological sort and how to compute it.



Given a graph $G = (V, E)$ we want to find a sorting for V (i.e. a map $S : V \rightarrow \mathbb{N}$) such that if $v \rightsquigarrow_G w$ (w is reachable from v) then $S(v) \leq S(w)$.

For example, topological sorts of this graph is (5 4 2 3 1 0) or (4 5 2 3 1 0).

Notice that:

- Topological Sorting for a graph is not possible if the graph is not a **Directed Acyclic Graph**
- There can be more than one topological sorting for a graph

Def. The **collapsed graph** of a graph $G = (V, E)$ is the graph $G_{SCC} = (SCCs, \{E(SCC_1, SCC_2) : \exists v \in SCC_1, \exists w \in SCC_2 (v, w) \in E\})$.

DFS has an important feature which is the final processing time f , and we can use it because the largest f corresponds to the lowest processing index in the topological sort. So we can use DFS and inherit the sorting induced by f . This is true only for acyclic graphs and takes $O(|V| + |E|)$.

Prop. $v \rightsquigarrow_G w \Leftrightarrow f[v] \geq f[w]$, where f is computed by DFS.

In the case of a cyclic directed graph you can compute Tarjan with more than one node.

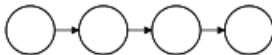
Transitive closure problem

→ [fischer_meyer_trans_closure.pdf](#)

→ [farzan_trans_closure.pdf](#)

Definition of transitive closure of graph. Reduction of the transitive closure problem to Boolean matrix multiplications. Fischer's and Meyer's algorithm: intuition, code, and complexity.

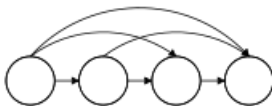
Input



Given $G = (V, E)$ we want to compute $G^* = (V, E^*)$, where $E^* = \{(v, w) | v \rightsquigarrow_G w\}$.

We want to model it almost as a matrix multiplication problem using the adjacency matrix.

Output



Transitive closure

We will see how the problem of finding the transitive closure of a graph is equivalent to the problem of boolean multiplication.

22/03/18

A naive algorithm is to compute $v \rightsquigarrow_G w$ for each $v, w \in V$, which has a very bad cost of $O(|V|^2) \cdot O(|V| + |E|) = O(|V|^3 \cdot |E|)$.

Fischer and Meyer algorithm

(for transitive closure problem)

1. compute the SCC collapsed graph, and go from a general graph to an acyclic graph
2. sort the nodes of the new graph in a way in which $AdjM(G)$ is upper triangular, and this can be done by using topological sort.
We know that if $S(v) \leq S(w)$ then w is not reachable from v , thus $AdjM(G)$ is upper triangular.
3. Now we can split it in 4 blocks $M = AdjM(G) = \begin{pmatrix} A & C \\ 0 & B \end{pmatrix}$ so we get $M^* = AdjM(G^*) = M \cdot \dots \cdot M = \begin{pmatrix} A^* & A^*CB^* \\ 0 & B^* \end{pmatrix}$
by using the transitive closure, not matrix multiplication.
4. Rebuild the transitive closure of G . This is easy because all the elements in a SCC have the same transitive reachability set.

Let's find the complexity of this algorithm:

1. is tarjan complexity $\Theta(|V| + |E|)$
2. is topological sort $O(|V|^2)$
- 3.
4. $O(|V|^2)$

$T(n)$ is the complexity of computing M^* where n is the size of the square matrix M , so

$$T(n) = \begin{cases} \Theta(1) & n = 1 \\ 2T(\frac{n}{2}) + O(n^{\log_2 7}) & n > 1 \end{cases}$$

by using master theorem we get $T(n) \in O(n^{\log_2 7})$.

The overall complexity of Fischer-Meyer is $O(|V|^2) + O(n^{\log_2 7}) = O(|V|^{\log_2 7})$.

This is not the best algorithm possible to perform transitive closure, but probably is this the best for very large and dense graphs. In the case of small graphs it is better to avoid the calculation of $AdjM$.

```

1 collapse(G, MN)
2 # MN is an array of arrays, each of them represents a new node
3 V2M <- zero(|V|)
4 i <- 0
5 # looping on all the nodes included into one single array
6 while i < |MN| do
7   for v in MN[i] do
8     V < MN[v] <- i
9     AdjR[i] <- []
10    i <- i+1
11  for v in range(1, |AdjG|) do
12    for w in Adj(v) do
13      AdjR[V2M[v]].append(V2M[w])
14  return AdjR

```

Adj list to adj matrix

```

1 Adj2AdjM(Adj)
2 # get the corresponding adj matrix
3 M <- zero(|Adj|, |Adj|)
4 for v in range(1, |Adj|) do
5   for w in Adj[v] do
6     M[v,w] <- 1
7  return M

```

Upper triangular matrix transitivity closure:

```

1 UTMATRIXTC(M)
2 # M is a square matrix
3 n <- n_of_rows(M)
4 if (log2 n ∉ ℕ) then
5   M <- fix_it(M)
6  return UTMATRIXTC_real(M)

```

```

1 UTMATRIXTC_real(M)
2 n <- n_of_rows(M)
3 if (n==1) then
4   return M
5 (A, B, C) <- split_three_blocks(M,n)
6 A* <- UTMATRIXTC_real(A)
7 B* <- UTMATRIXTC_real(B)
8 D <- MatrixMultiplication(A*,C)

```

```

9   D <- MatrixMultiplication(D, B*)
10  # replace elements  $\neq 0$  by 1
11  D <- zero_one(D)
12  return split_three_blocks(A*, B*, D)

```

```

1  Fischer_Meyer(G)
2   SCCs <- TarjanSCCs(G)
3   reverse(SCCs)
4   AdjSCC <- collapse(G, SCCs)
5   M <- Adj2AdjM(AdjSCC)
6    $M^*$  <- UTMatrixTC(M)
7   return decollapse_AdjM( $M^*$ , SCCs)

```

Notice that we have not implemented `decollapse_AdjM` yet.

→ [leggere qui](#) per capire transitive clousure

26/03/18

Weighted graphs

Definition of weighted graph and examples. Definition of the single source shorted path problem. Definition of (min/max)heap.

Definition of binary heap; complexity of building a binary heap, inserting a new key, extracting the minimum, and decreasing a value.

Fibonacci heap: complexity of building a Fibonacci heap, inserting a new key, extracting the minimum, and decreasing a value.

Dijkstra's algorithm: intuition, code, limitations, and complexity.

→ p.595-601 Cormen

graphs + a weighted function $W : E \rightarrow \mathbb{R}$ labeling the edges

Single source shortest path problem

Given a weighted graph (G, W) and $s, d \in G$ we want to find a path e_0, \dots, e_n from s to d such that:

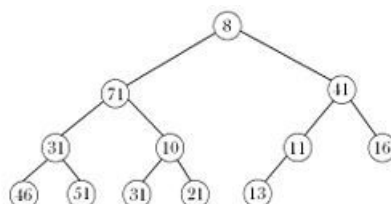
1. $e_0 = (s, v)$ for some $v \in G$ (s is the source)
2. $e_i = (v_i, v_{i+1})$ $e_{i+1} = (v_{i+1}, v_{i+2})$ for any $i \in [0, n-1]$
3. $e_n = (w, d)$ for some $w \in G$
4. $\sum_{i=0}^n w(e_i) = \min \sum_{i=0}^m w(\bar{e}_i)$ among all paths $\bar{e}_0, \dots, \bar{e}_m$ satisfying 1., 2. and 3. (the path minimizes the sum of the distances from the initial node)

Minheap is an abstract data structure that provides fast:

1. insertion
2. min estraction
3. key decreasing
4. deletion

Minheaps are useful because they support *logarithmic* insertion and deletion and can be built in *linear* time. The i -th location in the array will correspond to a node located on level $\lfloor \log i \rfloor$ in the heap.

The most famous implementation is a **binary heap**, which is an almost complete tree.



Definition:

- A min-max heap is a *complete binary tree* containing alternating *min* (or *even*) and *max* (or *odd*) *levels*. Even levels are for example 0, 2, 4, etc, and odd levels are respectively 1, 3, 5, etc. We assume in the next points that the root element is at the first level, i.e., 0.
- Each node in a min-max heap has a data member (usually called *key*) whose value is used to determine the order of the node in the min-max heap.
- The *root* element is the *smallest* element in the min-max heap.
- One of the two elements in the second level, which is a max (or odd) level, is the greatest element in the min-max heap
- Let x be any node in a min-max heap.
 - If x is on a min (or even) level, then $x.key$ is the minimum key among all keys in the subtree with root x
 - If x is on a max (or odd) level, then $x.key$ is the maximum key among all keys in the subtree with root x .
- A node on a min (max) level is called a min (max) node.

Dijkstra's algorithm

(for sssp)

```
1 initialize_single_source(G, W, s)
2   for v in G // for every node in G
3     d[v] <- ∞ // the distance of v from s is infinity
4     π[v] <- null // the predecessor of v in "minimal path" from s is null
5   d[s] <- 0
6   π[s] <- s
```

```
1 Dijkstra(G,W,s)
2   (d,π) ← initialize_single_source(G,W,s)
3   // build a minheap containing all the nodes of G
4   H <- makeMinHeap(V, key = d) // O(|V|)
5   while |H| > 0 do // O(|V|)
6     // u is a node s.t. d[u] = min{d[v] : v in V}
7     u <- extractMin(H) // O(log2(|V|))
8     for w in Adj[u] do // Θ(|E|)
9       if d[w] > d[u] + W[(u,w)] // Θ(1)
10        decreaseKey(H,w,d[u] + W[(u,w)]) // O(log2(|V|))
11        d[w] <- d[u] + W[(u,w)] // Θ(1)
12        π[w] <- u // Θ(1)
13   return (d,π)
```

notice that:

- if we use a negative weight it converges but the result is wrong
- if s is disconnected from the graph it returns infinite distance

The cost of Dijkstra is:

- $O((|V| + |E|)\log_2|V|)$ using binary heap
- $O(|V|\log|V| + |E|)\Theta(1)$ using Fibonacci heap

05/04/18

Binary heap in c++

→ *binary_heap.cpp*

nel branch *binary_heap*

We proposed a complete C++ implementation for binary heaps. Both correctness and complexity of the implemented methods were investigated. With the goal of implementing Dijkstra's algorithm, we also detailed associative binary heaps that allow to map the values insertion order into corresponding heap nodes. The heapSort algorithm was implemented too.

complexity? depends. The worst case is $i=0$: visit a branch of the complete tree in $O(\log_2 n)$

09/04/18

continuazione del binary heap...

☐ scaricare il file da moodle

12 aprile

A* algorithm (for sssp)

→ p.100-107 Cormen

→ *minimum_cost_paths.pdf*

A* algorithm: main idea, pseudo-code, and complexity. Differences between Dijkstra's and A* algorithms: case studies.

A^* = Dijkstra + heuristic distance

It consists in extracting from minHeap those nodes that minimize $d + h$.

- h is the heuristic distance, which estimates the distance between two nodes in the graph
- t is the destination
- $d[v]$ is the current best distance between s and v
- $\bar{h}[v]$ is the guest best distance between t and v
- $\pi[v]$ is the predecessor of v in the best path from s

```
1  A*(G, W, h, s, t)
2  for v in V
3      d[v] <- ∞
4       $\bar{h}[v]$  <- ∞
5       $\pi[v]$  <- None
6      d[s] <- 0
7       $\bar{h}[s]$  <- h(s, w)
8       $\pi[s]$  <- s
9      H <- BuildHeap(V, key =  $\bar{h}$ )
10     while (|H| > 0)
11         z <- DeleteMin(H)
12         if z == w
13             return RebuildPath( $\pi$ , w)
14         for u in Adj[z]
15             if d[u] > d[z] + w[z][v]
16                  $\pi[u]$  <- z
17                 d[u] = d[z] + w[u][z]
18                 H.DecreaseKey(u, d[u] + h(u, w))
```

alternatively using the closed set

```
1  while (|H| > 0)
2      z <- DeleteMin(H)
3      closed <- closed U {z}
4      if z == w
5          return RebuildPath( $\pi$ , w)
```

```

6     for u in Adj[z]
7         if d[u] > d[z] + w[z][v]
8             if u ∈ closed
9                 throw "your h sucks"
10            d[u] = d[z] + w[u][z]
11            π[u] ← z
12            H.DecreaseKey(u, d[u]+h(u,w))

```

☐ controllare se ho scambiato u e v

The complexity is the same of Dijkstra.

Floyd-Warshall algorithm

(for APSP)

→ p.570-576 Cormen

→ *boolean_matrices.pdf*

Definition and naive solution based on Dijkstra's algorithm. Floyd-Warshall's algorithm: main idea, limitations, pseudo-code, correctness, and considerations about memory usage. Paths reconstruction.

It solves **All pairs shortest path problem**: it finds the **shortest paths** in a **weighted graph** with positive or negative edge weights. A single execution of the algorithm will find the lengths (summed weights) of the shortest paths between *all* pairs of vertices.

$D^{(0)}, \dots, D^{(|V|)}$ are path distances

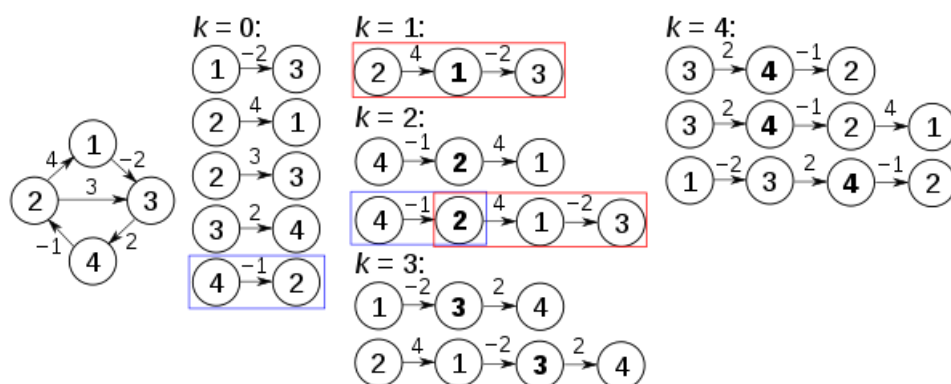
$$d_{ij}^k = \begin{cases} w_{ij} & \text{if } k = 0 \\ \min(d_{ij}^{(k-1)}, d_{ik}^{(k-1)} + d_{kj}^{(k-1)}) & \text{if } k > 0 \end{cases}$$

```

1 Floyd_Warshall(W)
2   D(0) ← W
3   for k in 1, ..., |V|
4     D(k) ← zeros(|V|, |V|)
5     for i in |V|
6       for j in |V|
7         D(k)[i][j] ← min(D(k-1)[i][j], D(k-1)[i][k] + D(k-1)[k][j])
8   return D(k)

```

the complexity is $\Theta(|V|^3)$



Floyd-Warshall example.svg

Routing problem and Contraction Hierarchies

→ [contraction_hierarchy.pdf](#)

Definition of the routing problem. Issues and opportunities on large graphs. The contraction hierarchies: motivations, basic idea, contraction algorithm, ordering heuristic, queries. Examples.

The method of **contraction hierarchies** is a technique to speed up [shortest-path routing](#) by first creating [precomputed](#) "contracted" versions of the [connection graph](#).

Contraction hierarchy non va implementato!

23 aprile

continuazione contraction hierarchy

Pattern matching

→ [p.923-931](#), 985 Corman

The pattern matching problem: definition. Naive algorithm: idea, pseudo-code, complexity, and examples.

The Knuth-Morris-Pratt's algorithm: the prefix(-suffix) function and its usage, pseudo-code, complexity, and working example.

Naive solution

Σ alphabet

Σ^* set of finite sequences of elements in Σ

$x \in \Sigma^*$ is a string

$|x|$ is the length of $x \in \Sigma^*$

if $x, y \in \Sigma^*$, $xy \in \Sigma^*$ is the concatenation of x and y

```
1 PM_Naive(T, P)
2   for s in [1, |T|-|P|] do
3     i <- 0
4     while i < |P| and P[i] == T[s+1] do
5       i <- i+1
6     if i == |P| then
7       yield s
```

$\Theta(|T| \cdot |P|)$

24 aprile

Knuth-Morris-Pratt

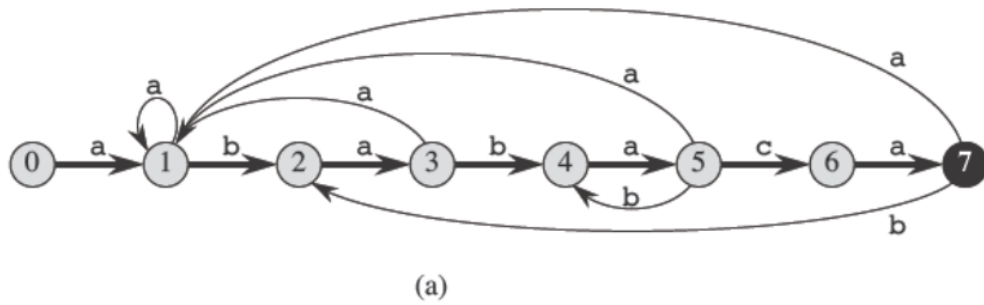
→ [pag.1002](#) Corman

→ [knut_morris_pratt.pdf](#)

string matching automata

$\delta : Q \times \Sigma \longrightarrow Q$

$\delta(q, a) = \sigma(P_q a)$



state	input			P
	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

(b)

i	1	2	3	4	5	6	7	8	9	10	11	
$T[i]$	a	b	a	b	a	b	a	c	a	b	a	
state $\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3

(c)

$\delta(i, T[i]):$

$$\delta(0, a) = \sigma(P_0 a) = \sigma(a) = 1$$

$$\delta(1, b) = \sigma(P_1 b) = \sigma(ab) = 2$$

$$\delta(2, a) = \sigma(P_2 a) = \sigma(aba) = 3$$

$$\delta(3, a) = \sigma(P_3 a) = \sigma(abaa) = 1$$

...

```

1 FiniteAutomationMatcher(T, δ, m)
2   n = length(T)
3   q = 0
4   for i = 1 to n
5     q = δ(q, T[i])
6     if q == m
7       pattern P occurs with shift i-m

```

its matching time on a text string of length n is $\Theta(n)$. This matching time, however, does not include the preprocessing time required to compute the transition function δ .

Transition function

$$\sigma: \Sigma^* \rightarrow Q$$

$$\sigma(x) = \max\{k : P_k \sqsubseteq x\}$$

```

1 ComputeTransitionFunction(P, Σ)
2   m = length(P)
3   for q = 0 to m
4     for a ∈ Σ
5       k = min(m+1, q+2)
6       repeat
7         k = k-1
8       until P_k ⊆ P_q a

```



```

9       $\delta(q, a) = k$ 
10     return  $\delta$ 

```

its running time is $O(m^3|\Sigma|)$

Knuth Morris Pratt

This algorithm avoids computing the transition function and its matching time is $\Theta(n)$, using just an auxiliary function π which we precompute from the pattern in time $\Theta(n)$ and store in an array $\pi[1, \dots, m]$

$O(|T| + |P|)$

26 aprile + 3 maggio

Boyer Moore Galil

The Boyer-Moore algorithm searches for occurrences of pattern P in text string T by performing explicit character comparisons at different alignments. Instead of a **brute-force search** of all alignments (of which there are $m - n + 1$), Boyer-Moore uses information gained by preprocessing P to skip as many alignments as possible.

Right to left matching

Characters in P and T are then compared starting at index n in P and k in T , moving backward. The strings are matched from the end of P to the start of P . The comparisons continue until either the beginning of P is reached (which means there is a match) or a mismatch occurs and one of the shift rules is applied.

Bad character rule

The bad-character rule considers the character in T at which the comparison process failed (assuming such a failure occurred). The next occurrence of that character to the left in P is found, and a shift which brings that occurrence in line with the mismatched occurrence in T is proposed. If the mismatched character does not occur to the left in P , a shift is proposed that moves the entirety of P past the point of mismatch.

```

- - - - X - - K - - -
A N P A N M A N A M -
- N N A A M A N - - -
- - - N N A A M A N -

```

```

1 def compute BCH(P)
2   for a in  $\Sigma$  do
3     BCH(a) <- 0
4   for i in [1, ..., |P|] do
5     BCH(P[i]) <- i
6   return BCH

```

Good suffix rule

Suppose for a given alignment of P and T , a substring t of T matches a suffix of P , but a mismatch occurs at the next comparison to the left. Then find, if it exists, the right-most copy t' of t in P such that t' is not a suffix of P and the character to the left of t' in P differs from the character to the left of t in P . Shift P to the right so that substring t' in P aligns with substring t in T .

If t' does not exist, then shift the left end of P past the left end of t in T by the least amount so that a prefix of the shifted pattern matches a suffix of t in T . If no such shift is possible, then shift P by n places to the right. If an occurrence of P is found, then shift P by the least amount so that a *proper* prefix of the shifted P matches a suffix of the occurrence of P in T . If no such shift is possible, then shift P by n places, that is, shift P past t .

```

- - - - X - - K - - - -
M A N P A N A M A N A P -
A N A M P N A M - - - -
- - - - A N A M P N A M -

```

Galil rule

The Galil rule is about exploiting periodicity in the pattern to reduce comparisons.

→ <https://stackoverflow.com/questions/38206841/boyer-moore-galil-rule>

```

1 def compute H(P,N)
2   k <- 0
3   for j in [1, ..., |P|-1]

```

```

4   if N(j) == j then
5       k <- j
6   H(|P|-j+1) <- k
7   return H

```

$N(i)$ is the longest suffix of $P[1..i]$ that is a suffix for P i.e. $P[i-j+1..i] = P[|P|-j+1..|P|]$
 $Z(i)$ is the length of the longest prefix of $S[i..|S|]$ that is also a prefix for S

```

1  def compute Z(P)
2      DZ[1], Z[2] compute explicitly
3      j <- 2, i <- 3
4      while i <= |S| do
5          if i >= j + Z(j) then
6              evaluate Z(i) explicitly
7          else
8              i' <- i - j + 1
9              Z(i) <- Z(i')
10             if Z(i) >= Z(j) - (i - j) then
11                 extend the tail of Z(i)
12             if j + Z(j) < i + Z(i) then
13                 j <- i
14             i <- i + 1
15     return Z

```

The complexity is sublinear in $|T|$: $O(|T| + \sum |P_i|)$

7 maggio

Longest common subsequence problem

→ p.390 Cormen

Notice that substrings do not have to be contiguous.

Naive solution

The complexity is $O(2^{|X|}|Y|)$ where $2^{|X|}$ is the number of substrings and $|Y|$ for the linear scanning of Y .

Recursive solution → p.393

8 maggio

DNA

String matching with finite automata

Suffix trie

http://www.cs.jhu.edu/~langmea/resources/lecture_notes/tries_and_suffix_tries.pdf
<https://www.hackerearth.com/practice/notes/trie-suffix-tree-suffix-array/>
<https://www.geeksforgeeks.org/pattern-searching-using-trie-suffixes/>

```

1  TOP <-  $\overline{t_1 \dots t_{i-1}}$ 
2   $\bar{R} \leftarrow \text{TOP}$ 
3  while  $\bar{R} \neq \perp$  do
4      if  $g(\bar{R}, t_i)$  is not defined
5          create  $\overline{Rt_i}$ 

```

```

6    $g(\bar{R}, t_i) \leftarrow \overline{Rt_i}$ 
7   if OLDR is defined
8      $f(\text{OLDR}) \leftarrow \overline{Rt_i}$ 
9    $\text{OLDR} \leftarrow Rt_i$ 
10   $\bar{R} \leftarrow f(\bar{R})$ 
11   $\text{TOP} \leftarrow g(\text{TOP}, t_i)$ 

```

```

1  def BuildTrie(T)
2    create  $\perp$ 
3    for  $t \in \Sigma$ 
4       $g(\perp, t) \leftarrow \bar{\epsilon}$ 
5       $f(\bar{\epsilon}) \leftarrow \perp$ 
6       $\text{TOP} \leftarrow \bar{\epsilon}$ 
7      for  $i \in [1, \dots, |T|]$  do
8         $\text{OLDR} \leftarrow \text{None}$ 
9         $\bar{R} \leftarrow \text{TOP}$ 
10       while  $g(\bar{R}, t_i)$  is None do
11         create  $R'$ 
12          $g(\bar{R}, t_i) \leftarrow R'$ 
13         if OLDR is not None
14            $f(\text{OLDR}) \leftarrow R'$ 
15          $\text{OLDR} \leftarrow R'$ 
16          $\bar{R} \leftarrow f(\bar{R})$ 
17       if OLDR' is not None
18          $f(\text{OLDR}') \leftarrow g(\bar{R}, t_i)$ 
19        $\text{TOP} \leftarrow g(\text{TOP}, t_i)$ 

```

- The BuildTrie procedure is optimal, i.e. it is $\Theta(|Q|)$ where Q is the set of nodes in $STrie(T)$
- $|Q| \in O(|T|^2)$

10-14 maggio (perse)

suffix trie pseudocodice

- ☐ chiedere gli appunti
- ☐ suffix tree <https://www.geeksforgeeks.org/pattern-searching-set-8-suffix-tree-introduction/>

15 maggio

Suffix array

pattern matching on a suffix array takes $O(|P|\log|T|)$, where P is the pattern and T is the string

```

1  def  $\leq l(T_1, T_2, l)$ 
2     $S_1 \leftarrow |T_1|$ 
3     $S_2 \leftarrow |T_2|$ 
4     $k \leftarrow 1$ 
5    while  $k \leq S_1$  and  $k \leq S_2$  and  $k \leq l$  do
6      if  $T_1[k] > T_2[k]$  then
7        return false
8      else if  $T_1[k] < T_2[k]$  then

```

```

9      return true
10     k <- k+1
11     return k = S1+1 or k = l+1

```

```

1  def suffixArrayMatching(T, SA, P)
2     $\phi \leftarrow |P|$ 
3    if  $\leq l(P, T_{SA[1]}, \phi)$  then
4       $L_P \leftarrow 0$ 
5    else if  $\leq l(T_{SA[|T|]}, P, \phi)$  then
6       $L_P \leftarrow |T| + 1$ 
7    else
8       $(L, R) \leftarrow (1, |T|)$ 
9      while  $R - L > 1$  do
10         $M \leftarrow (R+L)/2$ 
11        if  $\leq l(P, T_{SA[M]}, \phi)$  then
12           $R \leftarrow M$ 
13        else
14           $L \leftarrow M$ 
15       $L_P \leftarrow R$ 
16      compute  $R_P$ 
17      return  $(L_P, R_P)$ 

```

In order to improve complexity we may try to store the longest common prefixes among (P,L) and (P,R)

☐ studiare questo codice <https://www.geeksforgeeks.org/suffix-array-set-1-introduction/>

`txt+i` è semplicemente un puntatore all'i-esimo elemento dell'array txt, e mi permette di stampare il testo solo da un certo punto in poi

LCP array

the algorithm returns the LCP between T_1, T_2 and starts the comparison from char s

```

1  def LCP ( $T_1, i_1, T_2, i_2, s = 1$ )
2    while  $T_1[i_1 + s] = T_2[i_2 + s]$  and  $|T_1| \geq i_1 + s$  and  $|T_2| \geq i_2 + s$ 
3      s <- s+1
4    return s-1

```

pattern matching becomes:

```

1  def suffixArrayMatching(T, SA, P)
2     $\phi \leftarrow |P|$ 
3    if  $\leq l(P, T_{SA[1]}, \phi)$  then
4       $L_P \leftarrow 0$ 
5    else if  $\leq l(T_{SA[|T|]}, P, \phi)$  then
6       $L_P \leftarrow |T| + 1$ 
7    else
8       $(L, R) \leftarrow (1, |T|)$ 
9      l <- LCP(T, SA[L], P, 0)
10     r <- LCP(T, SA[R], P, 0)
11     while  $R - L > 1$  do
12       M <- (R+L)/2

```

```

13     h <- min(r,l)
14     if  $l \leq l(P, T_{SA[M]}, start = h, upto = \phi)$  then
15         R <- M
16         r <- LCP(T, SA[M], P, 0, h)
17     else
18         L <- M
19         l <- LCP(T, SA[M], P, 0, h)
20     LP <- L
21     compute RP
22     return (LP, RP)

```

This does not change the complexity of the algorithm

17 maggio

```

1 def LCP(T, i,  $\bar{T}$ , t)
2     k <- 0
3     while  $i + k \leq |T|$  &&  $\bar{i} + k \leq |\bar{T}|$  &&  $T[k + i] = \bar{T}[k + t]$  do
4         k <- k+1
5     return k-1

```

☐ studiare questo codice (copiato in locale): <https://www.geeksforgeeks.org/%C2%AD%C2%ADkasais-algorithm-for-construction-of-lcp-array-from-suffix-array/>

☐ scrivere una seconda versione adattando lo pseudocodice del prof e senza std library

? in c++ è il **conditional operator**: <https://www.geeksforgeeks.org/%C2%AD%C2%ADkasais-algorithm-for-construction-of-lcp-array-from-suffix-array/>

suffix array

... (forse ho perso l'inizio)

```

1 (L,R) <- (1, |T|)
2 while R > L+1 do
3     M <- (R+L)/2
4     if l>=r then
5         if  $L_{LCP}[M] \geq l$  then
6             m <- l + LCP(T, SA[M] + l, P, l)
7         else
8             m <- LLCP[M]
9     else
10        if  $R_{LCP}[M] \geq r$  then
11            m <- r + LCP(T, SA[M]+r, P, r)
12        else
13            m <- RLCP[M]
14        if m == P || P[m] <= T[SA[M]+m] then
15            (R,m) <- (M,m)
16        else
17            (L, l) <- (M, m)
18    return L

```

21 maggio

```

1 def UpdatePhase(T, SA, Prm, Bh, H)
2   // initialization
3   for each bucket (l,r) do
4     count[l] <- 0
5     for c ∈ [l,r] do
6       Prm[SA[c]] <- l
7   for each bucket (l,r) do
8     for d ∈ [l,r] do
9       count[Prm[d]] <- count[Prm[d]]
10      Prm[d] <- Prm[d] + count[Prm[d]]
11      Bh2[Prm[d]] <- true
12    for d ∈ [l,r] do
13      if Bh2[Prm[d]] then
14        l <- min(J:J>Prm[d], Bh[J] or not Bh2[J])
15        for f ∈ [Prm[d]+1, l-1]
16          Bh2[f] <- false
17    for i in [1, |T|] do
18      SA[Prm[i]] <- i
19      if Bh2[i] then
20        Bh[i] <- true
21    return (SA, Prm, Bh)

```

This takes time $\Theta(|T|)$ and after $\log_2 |T|$ applications of the procedure we get the correct suffix array, so building the suffix array takes time $\Theta(|T| \log |T|)$.

24 maggio

persa

28 maggio

Red black tree

→ p.308 Cormen

To delete a node with value v:

1. search for v in the tree (let M be the corresponding node);
2. Search for the immediate successor of v: i.e., search for the leftmost node (C) of the right child of M;
3. M.value ← C.value;
4. replace C with its right son, N;
5. if C and N, then color N black;
6. if N is not the root, then N has a parent, and there are some cases where a recoloring is needed (search for figures on textbooks);

```

1 def RB-delete(T,v)
2   z <- RB-search(T,v)
3   if z == none then
4     return
5   if z.right == leaf then
6     y <- z
7   else
8     y <- searchSuccessor(T,z)

```

```

9  y.rigth.parent <- y.parent
10 if y.parent == none then      # y was root
11     T.root == y.right
12 else
13     if is_left_son(y) then
14         y.parent.left <- y.right
15     else
16         y.parent.right <- y.right
17     z_value <- y_value
18 if y.right.color == red or y.color == red then
19     y.right.color <- black
20     delete y
21     return
22 RBT_delete_fixup(T,y.right)
23 delete y

```

```

1  def RBT_delete_fixup(T,N)
2  while(N ≠ T.root and N.color == black) do
3      if is_left_son(N) then
4          S <- N.parent.right
5          if S.color == red then
6              S.color == black
7              N.parent.color <- red
8              left_rotation(T, N.parent)
9              S <- N.parent.right
10         else # sibling is black
11             if S.right.color == black and S.left.color == black then
12                 S.color <- red
13                 N <- N.parent
14             else
15                 if S.left.color == red then
16                     S.left.color <- black
17                     S.color <- red
18                     right_rotation(T,S)
19                     S <- N.parent.right
20                     S.color <- N.parent.color
21                     N.parent.color <- black
22                     left_rotation(T, N.parent)
23                     N <- S
24             else # N is the right son
25                 # exactly as before but sides are switched
26                 N.color <- black

```

☐ studiare questo codice e implementarlo: [c red black tree insertion](#)

Esame

su appuntamento, preferibilmente in gruppo
two parts:

- homework, implement algorithms in c++

- per l'esame bisogna implementare tutti gli algoritmi fino a quelli sui grafi, e poi solo un altro a scelta di quelli dopo i grafi (tutti i codici dell'argomento)
- oral presentation about a research paper
 - proporlo prima al prof per email (circa 10 lucidi, 15/20 minuti)