

$$T(n) = 4 \cdot T(n/2) + O(n)$$

$$d = n^{\log_2 4} = n^2 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} O(n^2)$$

$$c = 1$$

$$(a_2 b_1 + a_1 b_2) \cdot 10^{m/2}$$

$$(a_1 + a_2)(b_1 + b_2) = a_1 b_1 + a_1 b_2 + a_2 b_1 + a_2 b_2$$

$$a \cdot b = \underbrace{a_2 b_2 + a_1 b_1}_{\text{BUT}} \cdot 10^m + \left((a_1 + a_2)(b_1 + b_2) - a_1 b_1 - a_2 b_2 \right) \cdot 10^{m/2}$$

BUT

are the same

So:

$$T(n) = 3 \cdot T(n/2) + O(n)$$

$$d = \log_2 3 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} O(n \cdot \log_2 3)$$

$$c = 1$$

are very small -> so it still takes time

-> In the Median String Problem -> we can compute the total hamming distance of the word => if the partial is already > than the best seen so far, continuing to explore useless.

B&B MEDIAN SEARCH

@ praticamente l'opposto.

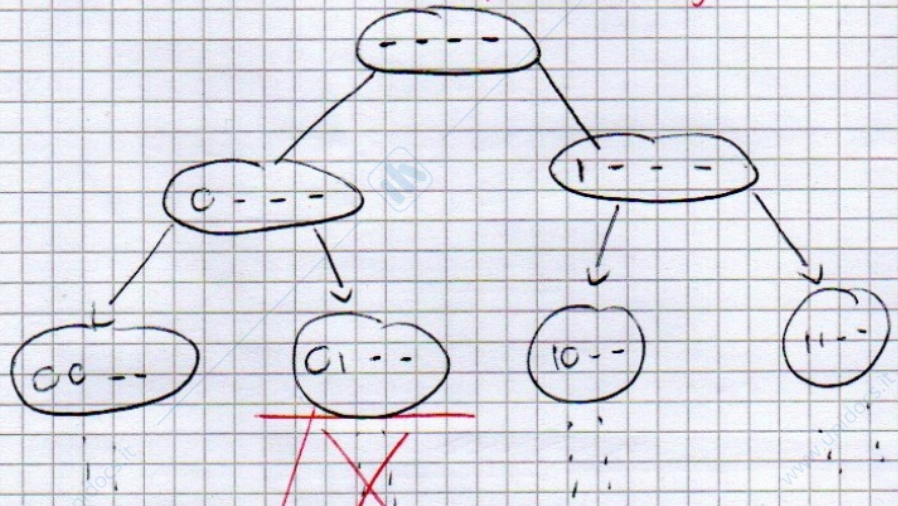
If the prefix distance > ~~the~~ best distance:
cut

else:

ofc here not sure if to find a better one but I cannot cut. I can try to see if I get a better hamming distance (shorter)

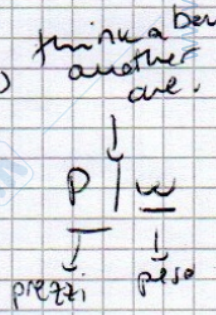
• Now we apply this to the knapsack problem: (just to give an intuition)

- $\{1, \dots, m\}$; $\{P_1, \dots, P_m\}$; $\{w_1, \dots, w_m\}$; T
objects ; ~~weights~~ *prices* ; *weights*



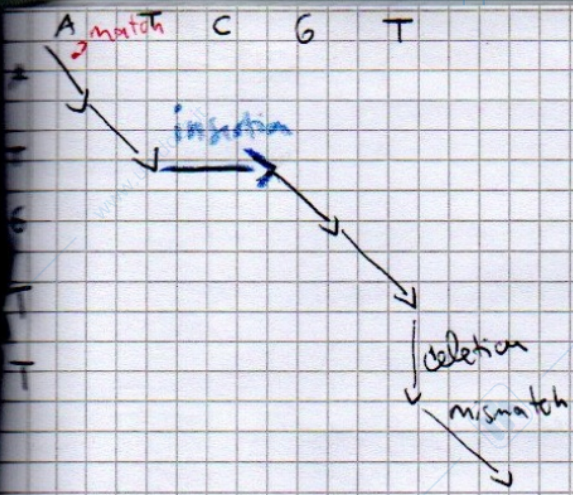
if already for the first choice we're exceeding the max weight
we are not interested in knowing how much it exceeds the max weight, we can simply cut.

-> one bound we have is P



-> If we're not looking for the best solution -> a good can work. it can become much faster.

Wester



con colori in slide molto chiari

it gives us info about the path through the grid

bc we're summarizing all the possible ways of aligning the first & 2nd string

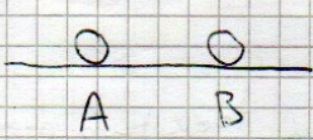
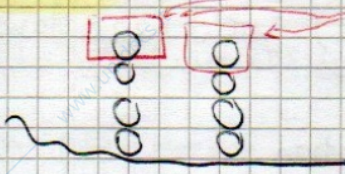
for establishing which alignment is better, we need a scoring function ← = computing the best path is equivalent to compute the best alignment in the grid.

A function that assigns scores to the paths

defining a scoring function is equivalent to assign weights to the path.

~~XXXXXXXXXX~~

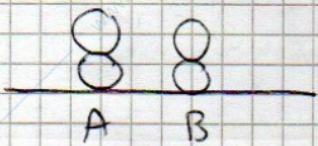
MOVES: take one rock or 2 rocks, ^{one of} each from each pile.



→ you need to play a strategy so that you'll be the last one to take the rocks.

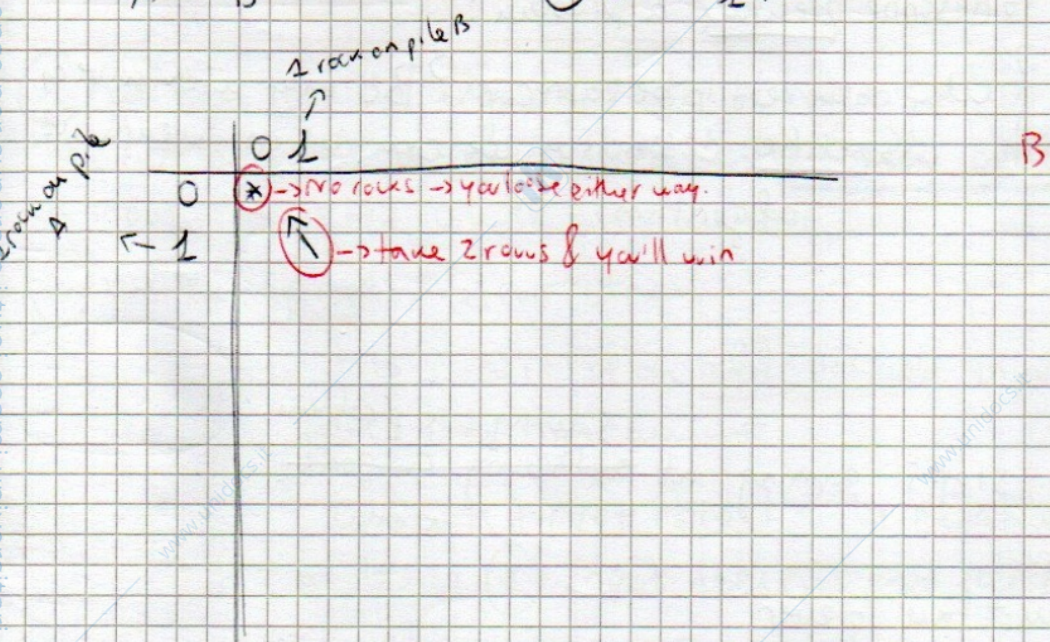
Let's say we have 10 rocks for each pile. If Alice starts, if we assume that ^{Bob} plays at his best, does Alice have a strategy for always being the winner?

Simple case:



- 2 options for Alice:
- ① takes 2 rocks → Bob wins
 - ② takes 1 rock → Alice wins

in reality ≠ inside a loop, it's 2+2 game wins for Bob.



The point of all this is that in simple cases we can come up w/ a situation, BUT if problems are more difficult we need a **strategy**.

↳ that's why we want to invent solutions for our problems.

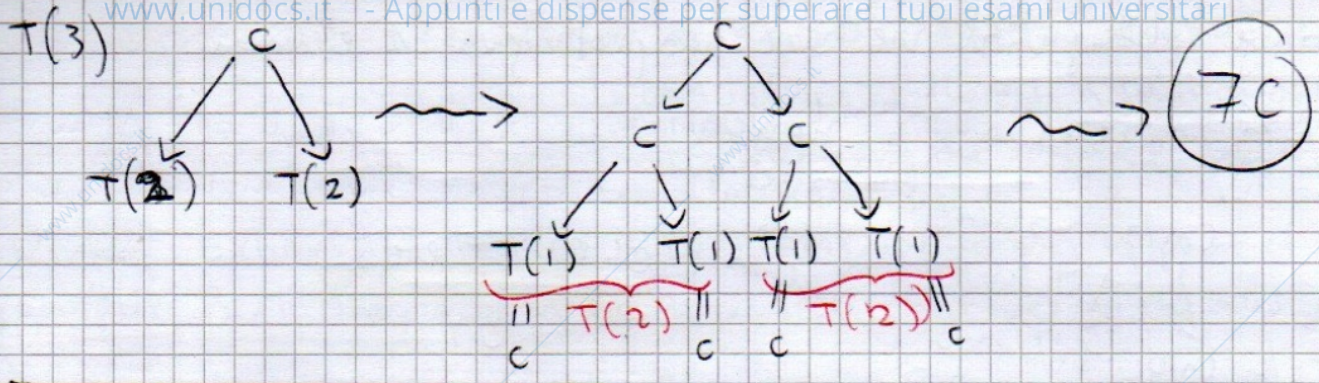
Combinatorial problem → def → slide.

ex. For us a **problem** is a set of pairs → made by input & output

ex. summing up 2 numbers:
 $(x, y) \rightarrow x$ is the input
 $\rightarrow y$ is the output.

for ex. $(\{2, 3\}, 5)$
 $(\{1, 4\}, 5)$

→



∇ $3c$ & $7c$ are just the n° of nodes of trees we have.

$2^n - 1 = \# \text{ NODES}$

Perhaps:

$T(n) \leq (2^n - 1) \cdot c$ ← **IS THIS TRUE? We verify by induction.**

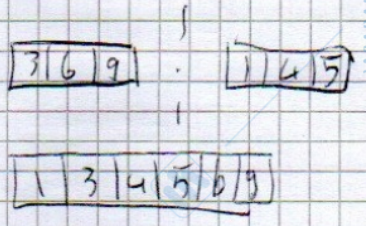
• **BASE CASE:** $n=1 \rightarrow$ we now that $T(1) = c$
 is $T(n) \leq (2^n - 1) \cdot c \rightarrow$ (Yes) $c \leq c$ ✓

• **INDUCTIVE STEP:** (guessing)
 $T(n) = 2 \cdot T(n-1) + c \leq 2 \cdot ((2^{n-1} - 1) \cdot c) + c$

$2^n \cdot c - 2c + c = 2^n \cdot c - c =$ ~~finisce copiare dei appunti da banca~~
 $(2^n - 1) \cdot c$ (OK)

MERGE-SORT(a)

if $\text{len}(a) = 1$: then | array of 1 el. \rightarrow return it.
 return

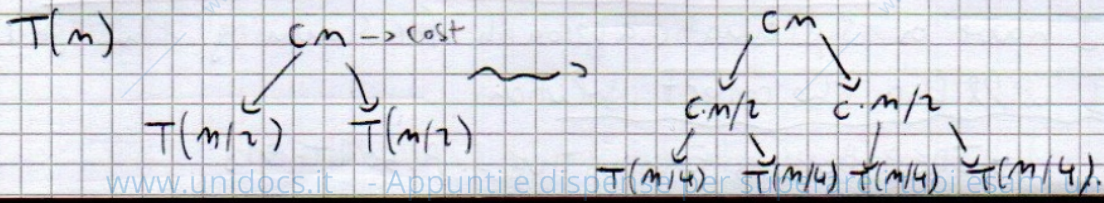


else:

MERGE-SORT(first half(a))
 MERGE-SORT(second half(a))
 MERGE(the 2 halves)

$T(1) = c$
 $T(n) = 2 \cdot T(n/2) + c \cdot n$
 for each element we will
 $2 \cdot T(n-1) \leftarrow$ Hanoi \rightarrow similar
 $2 \cdot T(n/2) \rightarrow$ this case.

\rightarrow does it maybe a difference



$f(n)$.

formal way to say it: $f(n)$ is $O(g(n))$ if there exists a point n_0 , after which $f(n) \leq c \cdot g(n)$ (i.e. there's a constant c , which multiplied by $g(n)$, will make $g(n)$ dominate $f(n)$).

EX.

$f(n) = 3n^2 + 4n \rightarrow O(n^2)$, bc we looked at "the most powerful" element which dominates. **BUT** it is also dominated by n^3 , bc $O(n^2) < O(n^3)$.

it is even more powerful $\Rightarrow O(2^n), O(n^3), O(n^2) \rightarrow$ all dominate $f(n)$.

BUT we're interested in the smallest upper bound.
we want the closest possible to $f(n)$

upper bound
they diverge from above

$O(n)$ is NOT an upper bound.

limit $f(n)$ from above

Big O notation is just a way for expressing this concept of upper bound.

! We could also bind functions from below.

$f(n) = \Omega(g(n))$
lower bound of $f(n)$

INTUITION: If n goes to infinity, $g(n)$ will be $< f(n)$.

DEF: $g(n)$ is a lower bound of $f(n)$ if there exists $n_0, s.t.$ that $g(n)$ scaled by c , will be lower than $f(n)$.

$f(n) = 3n^2 + 4n \rightarrow \Omega(n^2)$ & also $\Omega(n)$ but NOT $\Omega(n^3)$.

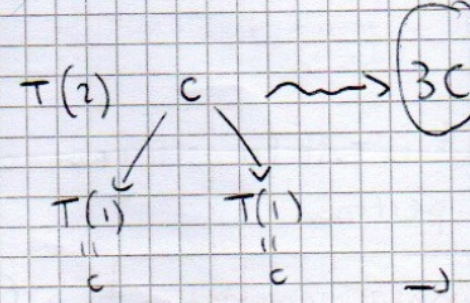
We can compare the 2 concepts:

$f(n)$ is $\Theta(g(n))$ if both $O(g(n))$ & $\Omega(g(n))$.
both a lower & upper bound

HANOI TOWER \rightarrow COMPLEXITY.

$T(1) = c$ (time for moving 1 disc) \rightarrow (we assume they're the same)

$T(n) = 2T(n-1) + c$



Solution at polynomial time will never be found.

BOOLEAN FORMULA:

$$(x_1 \vee x_2 \vee \neg x_3) \wedge (\neg x_3 \vee x_4 \vee x_5)$$

" x_1 or x_2 and not x_3 ..."

sol: $x_1 = x_2 = \text{true}$ $x_3 = x_4 = x_5 \rightarrow \text{false}$

$$(\text{true} \vee x_2 \vee \neg x_3) \wedge (\text{not false} \vee \text{false} \vee \text{false})$$

TRUE so enough so true \wedge true \checkmark

gli altri casi nelle slide

It is an NP-hard problem bc there's no clear path to follow in order to solve it.

0-1 KNAPSACK with integer weights \rightarrow NP-hard Problem

INPUT: m objects with integer weights

w_1, \dots, w_m and v_1, \dots, v_m values
band w on weight

OUTPUT: A subset $S \subseteq \{1, \dots, m\}$ s.t.

$$\sum_{i \in S} w_i \leq W \quad \sum_{i \in S} v_i \text{ is maximum}$$

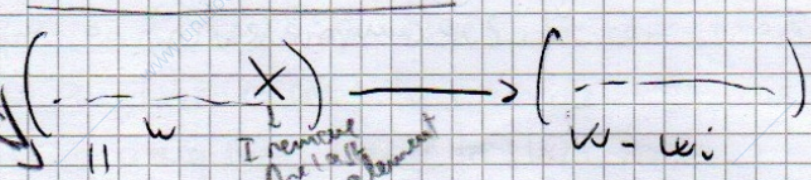
$OPT_{w,i}$: optimal value when we consider only the objects $\{1, \dots, i\}$ and their weight is exactly w .

(A) object i is part of the optimal solution

(B) object i is NOT part of the optimal solution

(A) $OPT_{w,i} \rightarrow OPT_{w,i-1}$ = specific weight that we consider

(B) $OPT_{w,i} = OPT_{w,i-1} \rightarrow OPT_{w,i}$ can be obtained from $OPT_{w,i-1}$



$OPT_{w,w} \rightarrow OPT_{w,i} - v_i$

ALGORITHM = a sequence of steps for solving a problem. It is still not a program. The program is the one implemented w/ a programming language. Basically also the pasta recipe is an algorithm.

↳ If many steps → humans are error prone. Instead a machine is faster & less errors prone.

↳ PROGRAM = description of an algorithm in a language, understandable by a computer

↳ *Calla differenza tra program & algorithm.*

• alphabet → a defined set of symbols

! If binary are strings: $\Sigma = \{0, 1\}$

ε 0101 0000 1 01101001 → same elements

"Sigma star" Σ^* indicates the set of all finite sequence of elements from Σ

• the edit distance $\delta(s, t)$ is the minimum number of changes necessary to turn s into t.

s: 0001010 t: 1000000

0001010 → 000010 → 000000 → 1000000
 (removed) (changed to 0) (added 1)

3 steps

$$\delta(0001010, 1000000) \leq 3$$

This only means that there's a way of transforming one into the other in 3 steps.

↳ at most 3! There could be other ways more efficient

• finite structures = domains w/ a limited (finite) number of symbols.

For this often these problems are called "combinatorials".

05/03/2021

The Rock Pile Game



• 2 piles of rock.

No such ~~box~~ ^{space}, you take another box & see if it fits.

FIRST-FIT (I)

$(U_1, \dots, U_m) \leftarrow (\emptyset, \dots, \emptyset)$ \rightarrow we take m boxes & initialize them as empty.

for $i = 1$ to n for each object $O(i)$ \rightarrow the box J has space for $O(i)$
 $J \leftarrow$ is the smallest index such that U_J has space for i

$U_J.add(i)$

return the PREFIX of (U_1, \dots, U_m) of non-empty boxes.

The overall complexity of FIRST-FIT is $O(n^2)$.

For this algorithm we'll never have more than 1 - less than half full box, bc even if we had 2 less than half full boxes, their contents would be put in the very same box.

LEMMA: FF returns an allocation in which AT MOST one box is full for less than 50%. \rightarrow cool property considering that it does NOT waste space.

$$H = \{U_J \mid \sum_{i \in S_J} s_i > 0.5\}$$

allocation

set of boxes full for more than 50%.

$$\sum_{U_J \in H} \sum_{i \in U_J} s_i > \sum_{U_J \in H} 0.5 \geq (FF(I) - 1) \cdot 0.5 = \frac{FF(I) - 1}{2}$$

for all boxes in H \rightarrow for every object in each box we want to sum their size.

$FF(I) - 1$ \rightarrow n° of boxes returned by FIRST FIT.

$\frac{FF(I) - 1}{2}$ \rightarrow that are fit for more than 50%.

~~part of boxes that are not more than 50%~~

$$\sum_{U_J \in H} \sum_{i \in U_J} s_i \rightarrow \text{size of the objects contained in each box of } H.$$

$$\sum_i s_i \geq \sum_{U_J \in H} \sum_{i \in U_J} s_i$$

this may be greater than that.

$$OPT \geq \sum_i s_i > \frac{FF(I) - 1}{2}$$

optimal n° of

$$OPT > \frac{FF(I) - 1}{2} \rightarrow 2 \cdot OPT > FF(I) - 1 \rightarrow 2 \cdot OPT + 1 > FF(I) \rightarrow$$

$$2 \cdot OPT > FF(I) \rightarrow \frac{2}{FF(I)} > \frac{FF(I)}{OPT} \rightarrow \text{APPROXIMATION VALUE.}$$

slide 21 -> spiega la complessità.

-> Rule of thumbs to evaluate recurrences -> NOT EASY!!!

let's assume we have a divide-and-conquer algorithm.

$$T(n) = a \cdot T(n/b) + f(n)$$

recurrence

a recursive call's

performs all on size n/b

↳ for reconstructing the final result.

$O(n \log_b a)$

$O(n^c)$

INTUITION: these 2 pieces contribute for a piece of the complexity; if $O(n \log_b a) > O(n^c)$, then complexity is $O(n \log_b a)$, otherwise & viceversa $O(n^c)$.
If comparable, a combination of the two. $d = \log_b a$

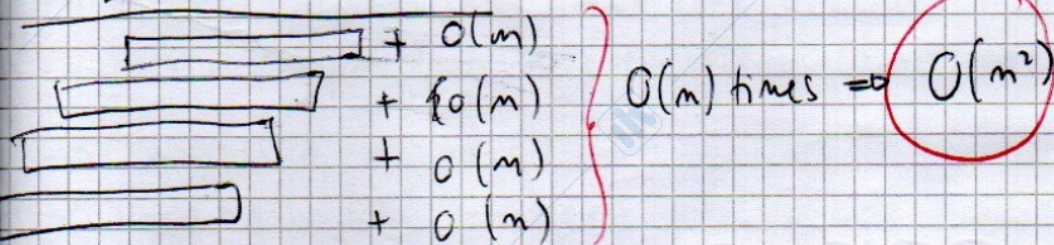
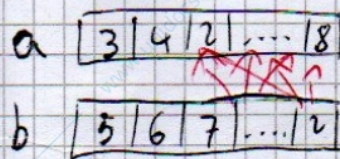
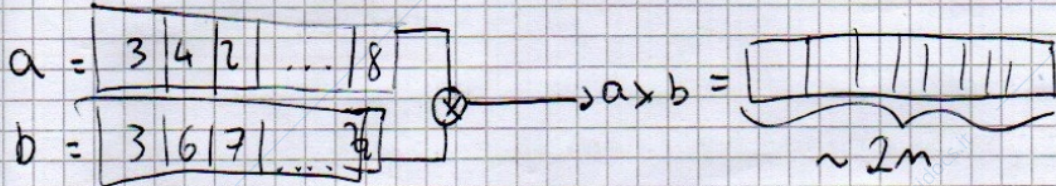
-> if $f(n)$ is $O(n^c)$ AND $c < d$, $T(n) = O(n^d)$

-> if $f(n)$ is $\Omega(n^c)$ AND $c > d$, $T(n) = O(n^c)$

-> if $f(n)$ is $\Theta(n^d)$, $T(n) = O(n^d \cdot \log n)$

Problem:

$a, b \rightarrow a \times b$



Now:

$$a = \begin{array}{|c|c|} \hline a_1 & a_2 \\ \hline \end{array}$$

$$a = a_2 + a_1 \cdot 10^{m/2}$$

$$b = \begin{array}{|c|c|} \hline b_1 & b_2 \\ \hline \end{array}$$

$$a \cdot b = (a_2 + a_1 \cdot 10^{m/2})(b_2 + b_1 \cdot 10^{m/2}) = a_2 \cdot b_2 + a_2 \cdot b_1 \cdot 10^{m/2} + a_1 \cdot b_2 \cdot 10^{m/2} + a_1 \cdot b_1 \cdot 10^m$$

↳ multipli catione

->

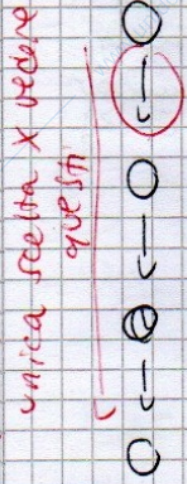
• we create a (grid), where edges have a w indicating how many attractions we can see from there (= by crossing that street)

→ (0,0) the starting vertex → source vertex (bc we cannot go upper no more to left)

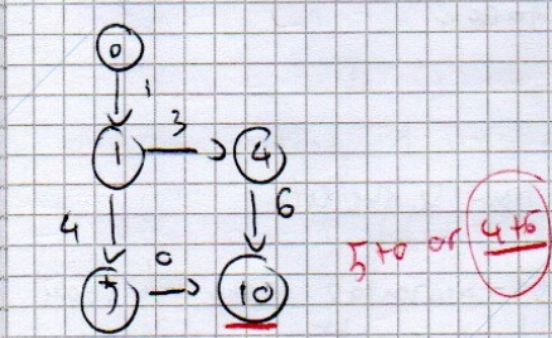
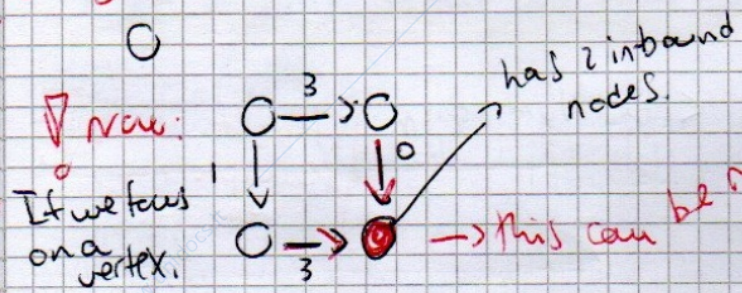
→ the last → the sink / target vertex (n,m)

• One drawback is that choices are done, if we're following a greedy approach, we would simply choose at each target the edge with highest weight, but it could be that this prevents us from seeing many attractions down to something

IDEA: we could build the best path by building best solutions for smaller parts of the grid
stessa cosa in cribi



• So we can get the best path from source to sink for ~~the~~ first column & first row.



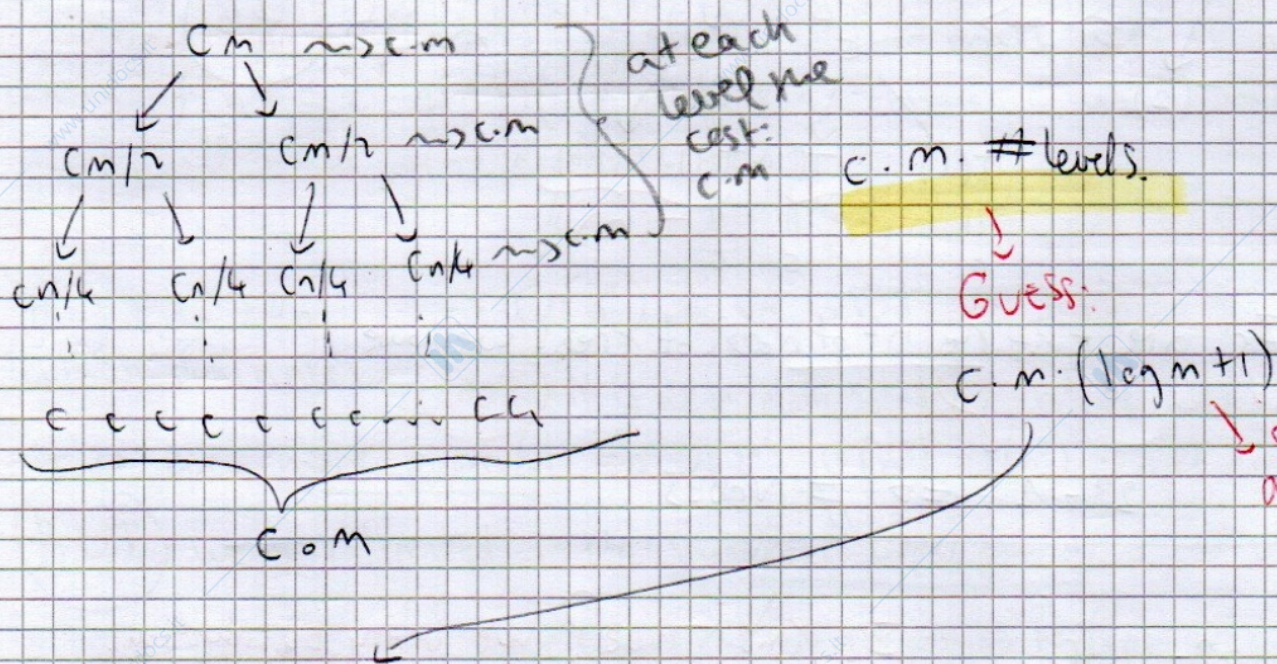
idea: let's focus on smaller instances & not directly from the first to the last.

↓
PROGRESSIVE APPROACH

→ input for algorithm: $\begin{matrix} 1 \\ \Delta \text{ matrix} \end{matrix}$ weights of the edges (horizontal) \vec{w}
 $\begin{matrix} 2 \\ 3 \\ 4 \\ 5 \end{matrix}$ $\begin{matrix} \uparrow \\ \downarrow \end{matrix}$ (vertical) \vec{w}

$\begin{matrix} 2 \\ 3 \end{matrix}$ → weights ↑
 $\begin{matrix} 4 \\ 5 \end{matrix}$ → weights →

we continue splitting until we have arrays of just 1 element



PROOF:

$$T(m) = c \cdot m \cdot \log m + cm$$

• BASE CASE:

$$T(1) = c \leq c \cdot m \cdot \log m + cm = c \cdot m = c \rightarrow \text{yes}$$

• INDUCTIVE:

$$T(m) = cm + 2T(m/2) = cm + 2 \left(\frac{cm}{2} \cdot \log \frac{m}{2} + \frac{cm}{2} \right)$$

$$cm + cm \log \frac{m}{2} + cm = cm + cm (\log m - \log 2) + cm =$$

$$cm + cm \log m - cm + cm = cm + cm \cdot \log m$$

$$O(m \cdot \log(m))$$

much smaller than $O(m^2)$

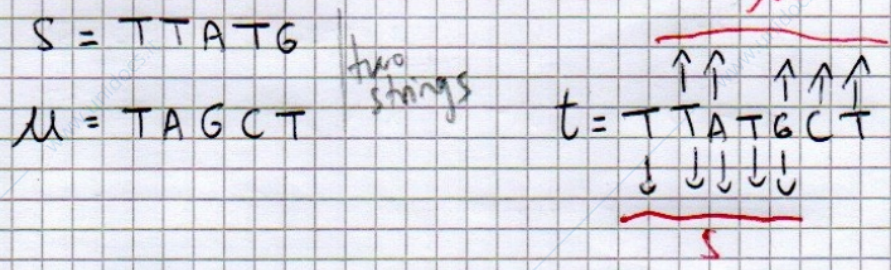
EXHAUSTIVE SEARCH ALGORITHMS

15/03/2021

- Some problems in which we cannot avoid to look at all possible search space; exhaustive search algorithms are the best candidates for these types of problems.

↳ I don't know the solution but I do know that the solution is in the search space, if I have a systematic approach for scanning this space sooner or later, I will find the correct solution.

↳ IDEA BEHIND EXHAUSTIVE S. ALGORITHMS.



the problem of largest common supersequence would be done by just combining the 2 strings; so it is kinda trivial.
 i.e. attaching them.

- INPUT: s, u strings
- OUTPUT: SCS \rightarrow a string t that is the shortest common supersequence

s' | prefixes of s & u we know that S = s₁ ... s_{len(s)}
 u' | prefixes of s & u u = u₁ ... u_{len(u)}

t' SCS of s' & u'

\Rightarrow a is the last character of t' \Rightarrow 3 possible cases

- (A) the character a is the last character of s' & u'
- (B) the character a is the last character of s' but NOT of u'
- ~~(C)~~ s''a = s' and u''b = u b \neq a
- (C) the character a is the last character of u' but not of s'
 s''b = s' and u''a = u b \neq a

if t' = t''a

- (A) in this case SCS(s'', u'') = t''
- (B) " " SCS(s'', u') = t'' \rightarrow bc the last character of t' a is coming from s'.
- (C) " " SCS(s', u'') = t'' \rightarrow

we have a recursion of optimality.

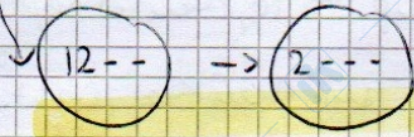
for each of these cases we can characterize the optimal solution to the smaller cases: to produce t', we simply need to take the best among these cases to start from.

THE IDEA: we have seen that an optimal solution allows us to find

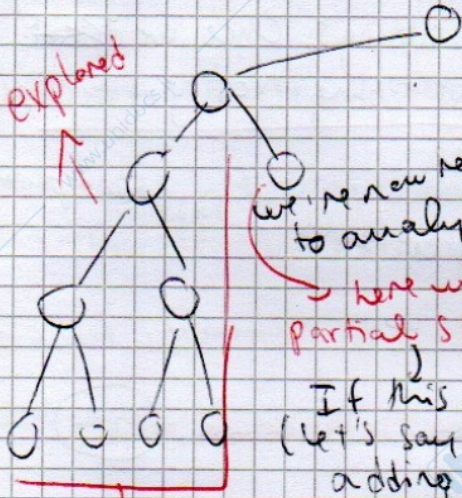
-> Uninteresting subtrees -> slide spiega bene

• For example let's assume we want to find bypass and go directly from node 10 to 17 -> we can do as before, to a node to a sibling of ancestor.

↳ **BUT** we need to know that exploring nodes 11-16 is useless.



• **BYPASS** -> simply as the end part in next vertex but ^{not} applied to a leaf but a node in the middle.



we're now ready to analyze this
here we set a partial score

If this is very low (let's say 10 & even by adding the best estimation score, the best alignment w/ what's below) & it is < 100, we can cut & go to the sibling of its ancestor

if we have an high score (given by the starting position)
ex. 100

$$\text{Score}(s, i, \text{DNA}) + (t-1) \cdot \underbrace{\ell}_{\text{hit for everything}} < \text{best score}$$

↑
length

then we can safely avoid to explore that node

BRANCH AND BOUND MOTIF SEARCH

6 we perform an optimistic score
7 if optimistic < best score
 Bypass
else
 Next vertex

we cut it, we don't ~~need to explore it~~

Even in the best case, we're not able to do better.

• This is a branch & bound technique

L_s = the optimistic score

HEURISTICS

• However ~~this~~ this is an improvement BUT NOT amazing -> sometimes cut

Section } Conclusion }

In section 2 we have seen how to write some mathematical formulas. Here we conclude that, blah blah dots
↳ displayed: "...".

The n² of sections is MANAGED BY LaTeX.

17/03/2021

optional -> assignment ①

↳ using Insertion Sort -> bc good for small lists.

↳ the function merge two lists. If all the elements are < than the ones of the other, it is not necessary to do it linearly, but you can do it in constant time by ~~appending~~ ^{appending} one to the other

-> The fact that the ^{most} gains depends on ① how the code is implemented in Alessandro's code in InsertSort ② the hardware is between 20-70.

APPROACHES : -> for measuring Performances.

Analytical

Experimental

- Asymptotic (it would tell us Merge Sort performs better BUT NOT the threshold).
↓
gives us the result for BIG inputs.
- Language & Machine Independent
- Worst case

Precise (bc we take timing of how long it takes the algorithm to perform)

-> we can analyze lengths of lists / inputs needed

e.g C++ has direct access to hardware factor.

Language & Machine Dependent

-> the language used the machine

Average case

-> Some algorithms, once fixed the length of the input, even if these is change perform always the same (ex. ~~Insert~~ Selection Sort)

Instead there are many algorithms that, even IF fixed the size of input, really depend on what the input is (i.e. is the input sorted? Then perform faster. An example of this is Insertion Sort).

↳ So we need to perform many tests, for understanding well how the algorithm performs.

- If the pairs are a finite number \rightarrow we can just list them.
- BUT in cases like the previous one (summing up two numbers) we'd need to have infinite combinations

\hookrightarrow we create a precise characterization.

we just specify a relationship, not how to obtain it.

How THE OUTPUT DEPENDS ON THE INPUT

the Role of the algorithm will be of finding an automatic way that executes ~~what does~~ & obtains the output

EX. The rock problem took two numbers as input.
(n, m)
 \hookrightarrow here the characterization is with the example of summing up 2 numbers.

Key concept: important to describe precisely a problem:

- what is the input
- what is the output

} why do we need to be so precise? Because only with a very specific description, we can design a capable algorithm.

ALGORITHM

\hookrightarrow Def. \rightarrow slide

\hookrightarrow Simply: sequence of steps aiming at ~~achieving a goal~~ ^{solving a problem}

finite \rightarrow it cannot take forever

unambiguous

- For describing algorithms we choose a language (unambiguous & precise enough)

we will use pseudocode \rightarrow it has a flavor of programming language but easier for us to understand.

- We want algorithms that always give us an answer.

DIFFERENCE w/ PYTHON:

for loops

for i = a to b
B

\rightarrow here a included & also b
 \downarrow
in Python a included & b excluded

would we need to arrive $b(\pi) = O(P) \cdot \underline{b(\pi)}$ steps. But in the worst case scenario, we're not able to reduce by 1 at every step \rightarrow you need to reverse the increasing to decreasing & then you can reduce by 1 \Rightarrow 2 steps each time.

So in the worst case scenario: $2 \times b(\pi)$.

\rightarrow Can the greedy approach be applied to

- we optimize only for the first 2 sequences, & then for all the remaining you follow your choice.
- we find the most similar only for the first 2 sequences.

no one knows the Approximation ratio

APPROXIMATION ALGORITHM \rightarrow takes less time BUT doesn't find the optimal solution.

ALGORITHM:

lines 3-4 \Rightarrow we're considering the first 2 DNA sequences

line 12 \Rightarrow once fixed the first 2, all the other DNA sequences are taken & looked for similarity with the match between 1 & 2.

\rightarrow there's then slide on the complexity. \rightarrow we pay efficiency w/ correctness

here concluded the part about this on the back. He now shows another algorithm

BIN PACKING PROBLEM \rightarrow very frequent.

we need to pack objects, each of which characterized by a ~~weight~~ size

- INPUT: $U = \{1, \dots, n\}$ objects

$i \Rightarrow s_i \in [0, 1]$ \rightarrow size for each object.

EX. $\{1, 2, 3\}$ $\{0.7; 0.4; 0.6\}$

\downarrow \downarrow
0.7 0.4
0.6

- OUTPUT: $\{U_1, \dots, U_k\}$ $U_i \in U$

\rightarrow we want to find a set of boxes such that:

$U_i \cap U_j = \emptyset \quad \forall i, j \rightarrow$ one object ^{only in} \checkmark 1 box, no more.

$\bigcup_i U_i = U \rightarrow$ all objects are placed in a box

$\sum_{i \in U_j} s_i \leq 1 \quad \forall j \rightarrow$ for each box no more than 1

- WORST case: for n objects, n boxes.

First fit algorithm:

• You put the object in the first box having enough space for it \rightarrow If there

This is correct if the smaller solutions are correct.

In order to show it, we must show it satisfies a PROPERTY via ~~induction~~ induction:

- 1 Identify the property \rightarrow the most difficult part, bc it may not be directly to algorithm correctness.
- 2 Show that it is satisfied in the base case
- 3 Assume that it is satisfied in all cases until ~~the~~ " $n-1$ "th one (inductive hypothesis)
- 4 Show that it is satisfied for the n th case, based on the correctness of the inductive hypothesis

ex. with hanoi tower \rightarrow property: the algorithm correctly moves the top most n disks from frompeg to topeg

This works easily for the base case but doesn't for the inductive case, it has too little preconceptions so we need a stronger property.

We add "if all disks in all pegs are correctly stacked and the top most n disks of frompeg are smaller than the top disks, any, in all other pegs"



fibonacci problem

input: an integer n

output: the n th fibonacci number $f_n = f_{n-1} + f_{n-2}$
considering $f_1 = f_2 = 1$

This can be solved through a recursion or in an iterative way

its correctness is simply proven, it uses only the definition. This version is less efficient, it would take more time.

to prove its correctness we need a property (called invariant bc it's true before the first iteration, during the for loop and after it). In this case an array is created so that in ~~propertion~~ position n of the array we always find the n th fibonacci number.

KNAPSACK PROBLEM → combinatorial. → fatto in programming

• we're a thief → we have a maximal amount of weight we can carry.
We want to steal the maximum worth.

↳ MAXIMIZING the max worth W .
max weight I can carry

- INPUT: $\{1, \dots, n\}$;
 m objects
 P_1, \dots, P_n ;
 ↓
 prices.
 W_1, \dots, W_n ;
 ↓
 weight.
 $T = \max W$

- OUTPUT: $X \subseteq \{1, \dots, n\}$ SUCH THAT
 ↓
 a set of objects

$$\sum_{i \in X} w_i \leq T \quad \text{and} \quad \sum_{i \in X} P_i \text{ is MAXIMUM}$$

the sum of weights of taken objects does NOT exceed T.

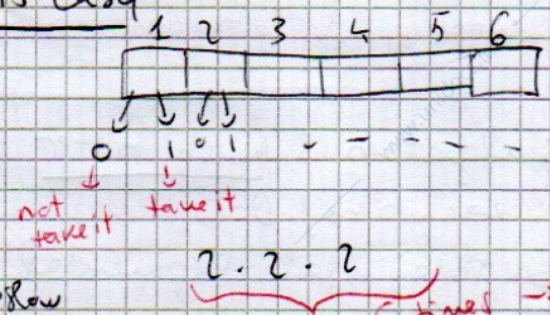
- (A) Provide an algorithm
- (B) Prove the algorithm complexity.
- (C) Evaluate its complexity.

A knapsack ($\bar{w}, \bar{p}, \bar{T}$)
 best value ← 0
 best_x ← {}
 for each subset $X \subseteq \{1, \dots, n\}$:
 ' if $\sum_{i \in X} w_i \leq T$:
 ' current_value ← $\sum_{i \in X} P_i$ (current value = sum of value)
 ' if (current_value > best_value):
 ' best_value ← current_value
 ' best_x ← X
 return best_x

B It is correct bc in the for we're looking at all possible possibilities, some point we'll find which is the best one.

↳ i.e. It is easy

C $\{1, \dots, n\}$
 for ↑ for merge cycle.
 $O(2^n \cdot c \cdot m)$



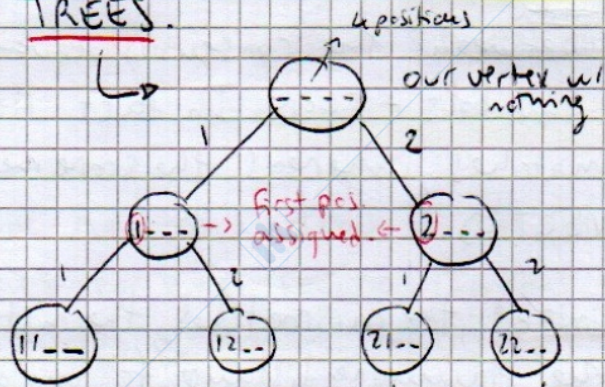
slow, when the input is big, it has to try many things.

$O(2^n)$ → EXPONENTIAL → slow → if $n=6$.

In this case we have the nucleotides (could be whatever symbol) that can be seen as tuples of numbers.

↳ To look in a TUPLE-ORGANIZED SEARCH SPACE

we use TREES.



and so until the 4 positions are occupied

We're interested in the last layer \Rightarrow the leaves of our tree.

-> Trees are very common in computer science.

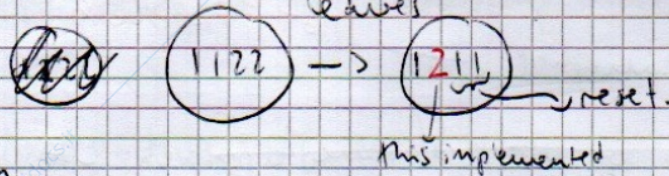
-> the **DEPTH** \rightarrow = how many layers we have \rightarrow can also be called "Height"

-> we're interested in trees with a fixed & constant n° of children (could be 2, 3, 4...)

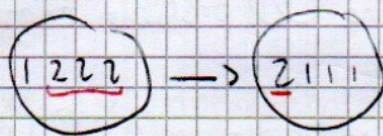
-> **BRANCHING FACTOR** = the n° of children departing from a node.

EXPLORING THE TUPLE SPACE

If we look at the ~~space~~ leaves it is like a counting:



or the same:



it is similar to ~~an~~ incrementing odometer.

We have a code, that given a leaf

Next leaf (a, L, k)

(in the previous ex.) $L=4, k=2$

- 1 for $i \leftarrow L$ to 1 we increment symbols from the last one to 1.
- 2 if $a_i < k$ ~~has this reached the last available~~
- 3 you implement if no
- 4 $a_i \leftarrow a_i + 1$
- 5 $a_i \leftarrow 1$ you put 1 if has reached the (k).
- 6 return a

ALL LEAVES \rightarrow just the idea that will be used in the motif finding problem

Although characterized in a \neq way, these 2 problems are computationally equal

-> OUR AIM now: finding relations in this 2 problems.

• We look for the Hamming distance between the consensus sequence w / the alignment matrix is all the symbols - score. In fact Hamming distance measures the n° of mismatches, whereas the score measures the matches. => HIGHLY CONNECTED.

↳ the consensus sequence minimizes the mismatches, the distance.

The consensus string minimizes the Hamming distance for an alignment matrix

-> We want to look for the string v minimizing the total distance with DNA \hookrightarrow OUTPUT of median string problem.

distance minimization
to Hamming distance
maximization to score

$$\min_u \min_s d_H(u, s) = \min_s \min_u d_H(u, s) = \min_s d_H(w_s, s) = \min_s (l \cdot t - \text{Score}(s, \text{DNA}))$$

Can be exchanged bc both minimization.

l.t is a constant so

we want to minimize a \ominus operation
 \downarrow
maximize

$$= l \cdot t - \max_s \text{Score}(s, \text{DNA})$$

↳ THIS IS EXACTLY THE MOTIF PROBLEM

apart from constants that we can ignore. \leftarrow

-> In the median string algorithm we use the same approach: we loop from AAA to TTT, so that we examine all possible combinations

↓
you optimize line by line \rightarrow therefore this can be done in $O(n \cdot l \cdot t \cdot 4^l)$

-> We've produced 2 Brute force algorithms, both of which are exponential. \rightarrow but in the last 4^l is the exponential & $4^l \ll n^t$

22/03/2021

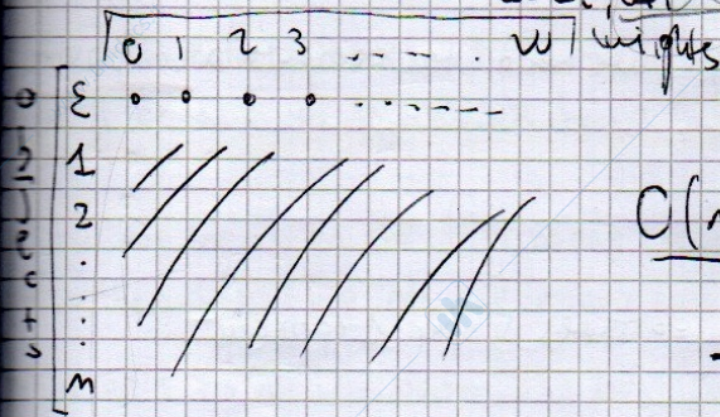
• Today we're going to see how to structure the algorithm, to avoid the exploration of useless paths in the search space.

\rightarrow we have to find an ORDER for looking at the search space.

TUPLES \rightarrow a more formal way for saying arrays.

\rightarrow

$$OPT_{w,i} = \max \begin{cases} OPT_{w,i-1} \\ OPT_{w-w_i, i-1} + v_i \end{cases}$$



$$O(m \cdot W)$$

THIS IS THE TRICKY POINT: the complexity of this looks polynomial but it considers the whole W .

size to represent W is

$$\log_2 W \text{ bits}$$

The number of bits is much smaller than the number W . So it could be

$$O(m \cdot 2^{\log_2 W})$$

the number of bits required for representing the W

It is said pseudopolynomial: polynomial in the value of the input BUT exponential in the size of the input. (required for representing it).

It is not polynomial, bc we need to consider the size, NOT the value of the input

DIVIDE & CONQUER ~ divide et impera

28/04/21

is a bit similar to the idea of dynamic programming.

We have seen an example of this, the mergesort algorithm.

The idea: if we're not able to compute the whole instance, we divide it.

MERGESORT(C) \rightarrow array \rightarrow we want to sort the elements of this

codice stile, chiaro

\rightarrow It is usually with the goal of gaining smth in terms of computation time or space usage.

1 $n \leftarrow$ size of C

2 base case \rightarrow if length $n=1$ I just return it, one element always ordered

The interesting part is list \rightarrow recombining of the lists, bc we are merging two

WE DO NOT VERIFY ALL OF THEM SIMULTANEOUSLY, but one by one. So having them all in memory is a waste of space, bc we don't need to access them all at one time. A way for having in memory at the time the one needed can allow me to save a lot of space.

↳ trick for exploring a huge search space.

we compute pieces of these data structures when needed.

keyword = yield

↳ ~~we~~ in our (ex.) (slide) we can assume yield means like return.

• it returns the value; **BUT** it remembers at which point of the computation we are.

For ex. for i in range generator(0, 4, 1):

○ dopo aver stampato 0, i è 1, quindi richiamo la funzione range generator e gli dico di iniziare da i a quel punto.

we resume the computation from where we left.

• nell'esempio di 100.000.000 → alla fine stampa così, ma the data structures are computed "on the fly" → quindi non sta creando una lista con 100.000.000 ma crea un elemento alla volta.

→ we can wrap generators w/ other generators ⇒ esempio range generator Filtered.

ASSIGNMENTS

19/04/21

↳ MARKING: 0-1-2-3-4

↳ ② assignments → da 0 a 4 → da 0 a 8 punti nell'esame.

8 voti → assignments

22 voti → scritto + orale

! sempre entrambi!!!!

LEGEND:

0 → missed submission / poor quality → se non mandi 1 file (pdf/code)

1 → solution has major issues - big mistakes / the quality of the report poor.

2 → more or less acceptable → solution can have some issues BUT ok report.

3 → solution has minor issues / good exposition, might lack some details. →

How to find a property for correctness? Not easy, the trick: try to abstract the idea that we can find it at first \Rightarrow continuous search.
 very tricky to prove smth by induction.

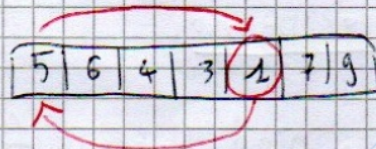
we're dealing with proving correctness

\rightarrow For iterative problems we make use of invariants i.e. some variables are kept TRUE

\downarrow
 we ~~keep~~ see an example.

SELECTION SORT(a, n) (we have an array, we scan it, looking for the smallest number, once found, we take it & put in the first position)

- 1 for $i \leftarrow 1$ to $n-1$
- 2 $a_j \leftarrow$ Smallest element among a_i, a_{i+1}, \dots, a_n .
- 3 Swap a_i and a_j
- 4 return a



can be "enriched" by `IndexOfMIN, a` function which helps in finding the minimum element.

INVARIANT property: once the initial part \rightarrow the elements of the first part are all sorted \rightarrow and \rightarrow smaller than those in the remaining part.

$$P(i) = (a[1:i] \text{ is sorted}) \wedge (a[1:i] < a[i+1:n]) \wedge$$

(a is reordering of the original array) \rightarrow tells us that no loss, just exchange

COMPLEXITY \rightarrow def. \rightarrow slide.

We need to do some assumptions (to simplify):

- Complexity is simply the n° of basic steps executed
- each step takes the same amount of time.

slide

\rightarrow It is just for us to have a rough idea of the time needed for that algorithm to perform.

Asymptotic Notation $\rightarrow f(n) = O(g(n))$

\rightarrow in this case what we mean is that when n grows a lot, $g(n)$ dominates

So any reversal $p(i, j)$ can decrease at most $b(\pi)$.

$$d(\pi) \geq \frac{b(\pi)}{2}$$

DISTANCE of the permutation π from the identity.

minimal computation necessary for deleting all breakpoints \rightarrow (if always best case).

ALGORITHM:

- 1
- 2 we consider all possible breakpoints & choose which reversal ensures to delete the most of breakpoints.
 - ASSUMPTION: we're always capable of getting rid of breakpoints.
 - ↳ it works only if we've the guarantee that at each step we can find a reversal reducing the n° of breakpoints.

IS IT POSSIBLE?

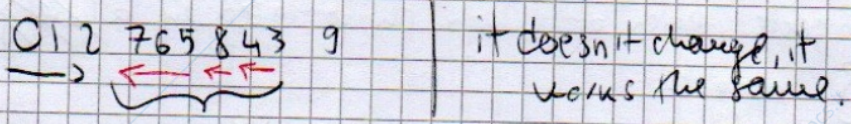
THEOREM \rightarrow slide

if we've a decreasing strip we're capable of decreasing at least by 1 the n° of breakpoints.

↳ PROOF: k is the smallest n° appearing in a decreasing strip, in our case $k=3$, therefore $k-1$ must be in an increasing strip (2).
 • $k-1$ must be at the end of the strip.

↳ By reversing, $b(k)$ will be close to $2(k-1)$.

It may be the case that 012 is on the right side & not a here:



IMPROVED BREAKPOINT REVERSAL SORT(π)

(line 5) in the else condition

in case we have only increasing strip BUT w/ breakpoints, what we do is to reverse an increasing, so that later we're sure we can find a reversal.

APPROX RATIO = 4 \rightarrow we can say that in the worst scenario every 2 steps we're capable of decreasing the n° of break points by 1.

$$\frac{2b(\pi)}{d(\pi)}$$

\rightarrow our solution.
 \downarrow
 optimal sol.

we record that: $d(\pi) \geq \frac{b(\pi)}{2}$

$$AR \leq \frac{2b(\pi)}{\frac{b(\pi)}{2}} = 4 \rightarrow \text{it is constant (the error)}$$

Bc $b(\pi) =$ n° of breakpoints present at the beginning, the possibilities are 2: there's a decreasing strip (we can find the reversal to reduce the n° of breakpoints by 1) or not.

If this were always the case, how many steps \rightarrow

if num coins + 1 < bestNum coins
I store it in a variable

Each time I need to recompute things however -> waste of time

IDEA: Since going from bigger to smaller, n° we have to recompute, we have to go from smaller to bigger ^{numbers} ~~times~~. have seen in the Fibonacci case

- ex) I want 1 -> 1 coin of 1
- I want 2 -> 2 coins of 1
- I want 3 -> I could have 3 coins of 1 BUT the best is 1 coin of 3

We use an array to memorize the best sol.
 So that we can have access to them & then everytime I need it I read it from memory.
 Trick: ~~adding~~ we're memorizing the optimal solution for the smaller cases

DPC change

- bestNum coins₀ ← 0
- per title le quantità di m.
- ~~...~~

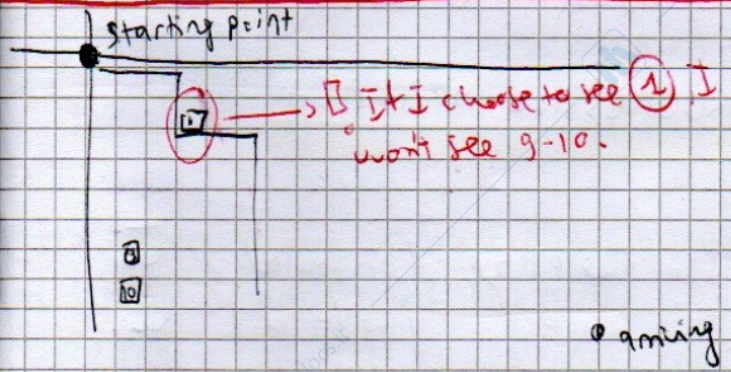
5. If m (the value considering now) > c_i:
 if bestNum coins_{m-c_i} + 1 < bestNum coins_m

if better of what seen up to now, we memorize it.

return bestNum coins_m.

PRICE: the memory -> you memorize what's done up to now

THE MANHATTAN TOURIST PROBLEM



we want to traverse the city

- Constraints:
- horizontal road only from left to right
 - vertical only from north to south.

-> During our travel, we want to maximize the n° of attractions that we see. We need to find the best path to ~~do~~ so.

We copy the subsequences in an **alignment matrix**: I look at it column by column \Rightarrow I basically count how often the nucleotides appear in each column.

\hookrightarrow "PROFILE MATRIX"

\rightarrow the **CONSENSUS STRING** is formed by the most frequent symbols column by column.
 \hookrightarrow we obtain the most similar to each of the substrings.

To measure how good this is we define a **SCORE**: the sum of the highest frequencies for each column. (In this case (42): it corresponds to the match it has with the alignment matrix).

the higher it is, the higher the n° of matches
& the higher it is, the lower the variance.

• this is a deterministic procedure.

once we pick the starting positions, this table is determined & so is the consensus sequence.

BRUTE FORCE MOTIF SEARCH (DNA, t, n, l)

1 best score $\leftarrow 0$ we define this variable

2 for each (s_1, \dots, s_t) from $(1, \dots, 1)$ to $(n-l+1, \dots, n-l+1) \rightarrow$ observe again

viewe poi spiegata la complexità $O(t \cdot l \cdot n^t)$
line 2
line 3

very slow
 \downarrow
So other perspective

other perspective: \rightarrow (we will reach a faster algorithm).

Median strings

\hookrightarrow we get the Hamming distance \rightarrow the n° of positions at which they differ (2 strings)

• We then define the **total distance** between a string u with a DNA matrix is the minimum distance that we can achieve by considering all the candidate starting positions

basically we're trying to find the string repeated more times w/ the smallest ~~number of~~ variations

Motif Finding Problem & Median String Problem are actually the same.
we want to find the positions w/ less variations in DNA. \hookrightarrow we want to find the string repeated w/ less variation.

→ up to now → we dealt w/ EXHAUSTIVE SEARCH, which are good bc they give us correct results, BUT pretty slow & with high complexity.

↓ In real scenarios, NOT much used, BUT can be improved

→ Greedy algorithms → each time has to perform a choice, they tend to check the solution looking the best at the moment.

Dynamic P.:

• IDEA: for building the optimal solution we partition the instances in sub-instances. We solve each of them in polynomial time

↳ ex. the change problem

→ If in input we have M , the best solution can be obtained from a number $n < M$? If yes, we can use dynamic programming approach.

↳ ex. $M=77$, if somebody tells us the best solution for 76, as there are no coins $=1$, we can obtain the optimal for 77.

↳ intuition: usually optimal solutions are kinda recursive; we can build solutions starting at from optimal solutions of smaller cases

ex. $c = (1, 3, 7)$
(possible coin sizes)

if (i_1, i_2, i_3) not optimal sol for 76 (we assume by contradiction) there would be a solution

$$j_1, j_2, j_3 < i_1 - 1 + i_2 + i_3$$

& by adding 1 we'd get:

$$j_1 + 1 + j_2 + j_3 < i_1 + i_2 + i_3$$

↳ contradicting i_1 (slide)

→ we will ask the best sol. for cases 76, 74 & 70 $(-1, -3, -7)$. We ask to someone else, we receive the best solutions for these 3 subcases. We know that the best sol. for 77 will be adding 1 coin to some of these 3, by taking the best one.

• What does it mean to ask to someone else? A recursive call, we essentially launch the algorithm to a smaller instance.

- In alto a destra -> new project
- scrivi a sinistra e vedi a destra cosa puoi vedere.

-> You describe in LaTeX how to write -> sort of programming languages.

```
\documentclass[a4paper]{article}
```

↓ deve sempre essere messo
 ↓ format
 ↳ type of writing
 ↳ mett sempre

```
\usepackage[utf8]{inputenc}
\usepackage[T1]{fontenc}
```

↳ non ci frega cosa è -> mett.

```
\usepackage[english]{babel}
```

↳ definisci la lingua d'uso.

↳ LaTeX is able to automatically cambiare a capo

```
\usepackage{a4wide}
```

↳ aggiunge un po' di spazi, viene meglio.

DEVE SAPERE LA LINGUA CHE SI USA

breaks the word
↳ See come Subdivide in syllabus.

```
\usepackage{modern}
\usepackage{authblk}
```

TUTTE QUESTE METTILE, YOU DON'T CARE WHAT THEY ARE.

```
\title{An example document in \LaTeX}
\author{name}
\author{name}
\affil{DISI, University of Bologna, Italy}
\date{10 March 2021} % \date{\today}
\begin{document}
```

-> see you use, aggiungi \date()

```
\maketitle
\end{document}
```

```
\section{Introduction}
This is an example of introduction.
```

```
\section{Body} \label{sec-body}
```

in all the sections puoi chiamarlo \ref{sec-body}

It is easy to write Mathematics here
For example: ↳ se lasci una empty va a capo, altrimenti dai solo spazio.

mostre:

$$\sum_{i=1}^n \frac{n(n+1)}{2}$$

```
\[ \sum_{i=1}^n i = \frac{n(n+1)}{2} \]
```

lists that are sorted. \rightarrow So it can be done in linear time ! If not ordered it would take us quadratic time.

\rightarrow To be sure the complexity gains smth, we need to compute a complexity analysis

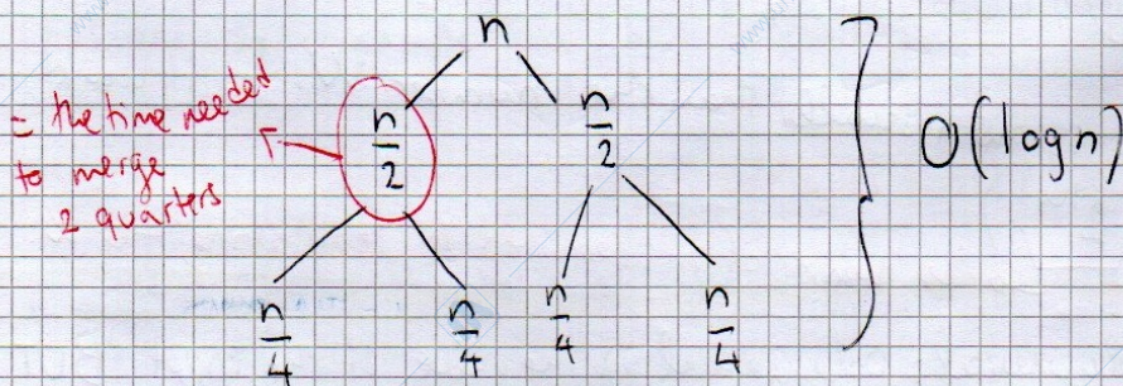
↓ For the recursion:

$T(1) = 1$ \rightarrow if size = 1 \Rightarrow takes simply 1 operation

$$T(n) = \underbrace{2T(n/2)}_{\text{OK}} + cn$$

\hookrightarrow time for merging the two.

\rightarrow A recursion tree is a tree depicting how execution is carried out.



\rightarrow The overall complexity for this is $O(n \log n)$

! Merge Sort is an optimal ~~problem~~ ^{solution}, bc there are mathematical proofs that show that there's no better complexity for ordering an array.

\rightarrow If we focus again on the global alignment problem: If we are interested in the n^2 , we need quadratic time bc we have to run but we can use linear space, bc only the last column is needed

\rightarrow If we want the actual alignment, can we do better than quadratic space? We need to exploit the trick that the score of the distance can be calculated in linear space.

slide: the complexity of PATH

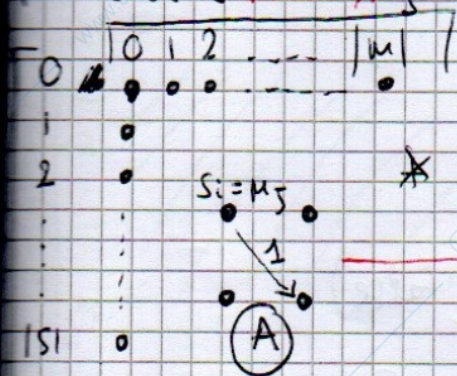
\rightarrow using a limited amount of space.

\rightarrow We want to compute the alignment from source to sink; any path has to go through the middle column at some point. So also the best path has to traverse it. However we don't know exactly when.

OUR AIM: try to compute where's the point in which the best path crosses

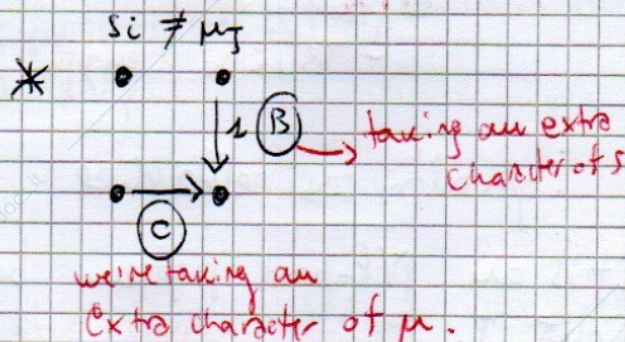
optimal solutions for the subcases.

→ As seen for previous problems, we may need a helper data structure, in this case a matrix



c_{ij} = length of the LCS of a i -prefix of s & a j -prefix of m .

ex. $c_{i,0} = i$
 $c_{0,j} = j$



$$c_{ij} = \min \begin{cases} c_{i-1,j-1} + 1 & \text{if } s_i = m_i \\ c_{i-1,j} + 1 & \text{if } s_i \neq m_j \\ c_{i,j-1} + 1 & \text{if } s_i \neq m_j \end{cases}$$

Complexity: $O(|S| \cdot |M|) \rightarrow$ polynomial.

ex. $s = TGC$
 $m = \epsilon$
↳ empty

needs to be a supersequence for s & m ; an empty one could be a superstring for m BUT NOT FOR S !

$t = TGC$

→ The difference with exhaustive search is that in dynamic programming we do reuse things → the work for computing sub-instances is reused. In exhaustive search scenario each time the optimality is checked from scratch.

→ The NP-hard problems → for them at the current state of art, the only known solution is at exponential time. Up to now no polynomial algorithm has been shown BUT still not proven that it is NOT possible to find a solution. (i.e. exhaustive search).

live in a limbo

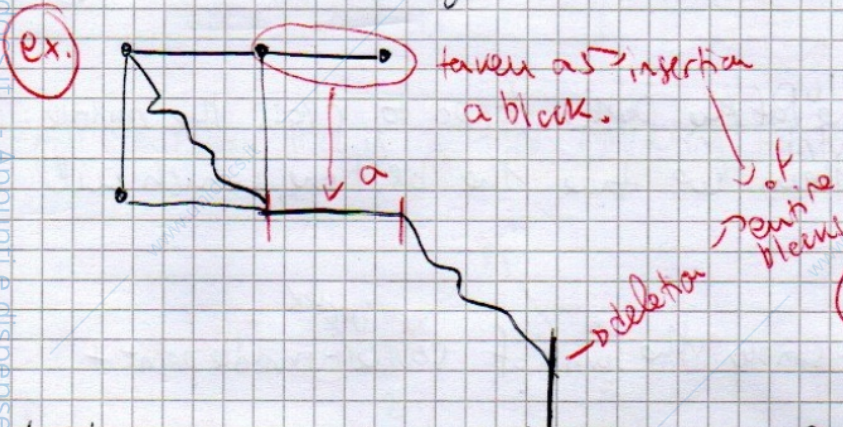
↳ Security, cryptography are based on the common assumed idea that no

→ for the global alignment we've seen that we can reduce the space complexity; what we want to do today is to find out if we're capable of reducing the time complexity through Divide-and-conquer.

↳ For the global alignment still researches about if we can reduce the $O(n^2)$ complexity to a linear one!

→ But we can look at a similar problem, the **BLOCK ALIGNMENT**

→ Let's consider a string $u = u_1 \dots u_n$, a t block string is so if n is a multiple of t.



• a string $u = u_1 \dots u_n$ is a t-block string if n is a multiple of t. So we can divide it in $\frac{n}{t}$ blocks of length t.

↳ So the length of the strings in this problem are under this constraint!

t-blocks are inserted or deleted as a whole, when you align them, they don't have to be the same. So:

$$S_{i,j} = \max \begin{cases} S_{i-1,j} - \alpha \text{ block del} \\ S_{i,j-1} - \alpha \text{ block insert} \\ S_{i-1,j-1} + \beta_{i,j} \text{ align} \end{cases}$$

↳ recursion which allows us to fill the dynamic programming data structure

these values should be computed before for all possible pairs of blocks / so that they are ready to be used in the matrix filling

→ A lookup table → where we store the values for aligning 2 blocks of length t. In this way we don't have to compute it everytime

! The values in the table must be stored in an alphabetical order, bc in this way we can perform binary search in the array (we don't use hash tables bc in the case of long & unbalanced hash clashes, we would need to perform a linear search in that list & we want to avoid this).

we put a strong bound in the worst case (in the worst case we take logarithmic time).
↳ (In the worst case, so not likely but we don't have the guarantee)

we guarantee that we remain in logarithmic time.

→ After this pre-computation we simply perform a change.

• **DECISION PROBLEM** → the solution of these is YES or NOT.
↳ can tell us if there's a solution or NOT.

• **OPTIMIZATION PROBLEM** → you look for the best, the optimal solution.
↳ from a

• Usually **EXHAUSTIVE SEARCH ALGORITHMS** are complex, point of view very complex, BUT they give us the guarantee that we find the answer.

↳ in real life → slow, so practically never used, BUT we're interested because it's an easy-idea method, which can constitute the base for further algorithms.

RESTRICTION MAPPING

↳ lowi fa esempio con motorways → in DNA same for finding restriction sites for restriction enzymes

Complete vs partial digest → we're NOT 100% sure that the enzymes cut in all restriction sites → usually they cut in ^{consecutive} ~~successive~~ points BUT we could obtain fragments with length 4, ignoring one site.

↳ what we obtain is cuts w/ ≠ lengths, not necessarily in 2 consecutive points. ↳ If length = 6 → this tells us that there are 2 points at distance = 6, BUT it may be that these are NOT consecutive → just a distance between 2 points (ex. 4-10).

• **MULTISET** = elements in a set can have multiplicity (ex. 2 has multiplicity 2, 3 also 2).
↳ (see in slide).

Partial Digest Problem

The 2nd biggest element 8 tells us that there are 2 points whose distance is 8. So it could be a point in 2 (2-10) or in 8 (0-8). THIS IS THE IDEA \rightarrow we don't know which of the 2 pts is the right one, we will now test it.

INCREMENTAL STRATEGY

$$L = \{2, 2, 3, 3, 4, 5, 6, 7, 8, 10\}$$

$$X = \{?\}$$

$$L = \{2, 2, 3, 3, 4, 5, 6, 7, 8\}$$

$$X = \{0, 10\}$$

I need to justify a distance of 10.

valid side \Rightarrow we encounter a decision point for 8 \rightarrow bc you can take it from left or right.

Then for 7 again:

from left:

$$X = \{0, 2, 3, 10\} \text{ but } 3-2=1$$

so $1 \in L$ not satisfied \rightarrow wrong, I take the other. (the right)

Then again the same for distance (6).

IDEA: By looking at the biggest distance & based on that I have 2 possible eligible options \rightarrow I test them.

Going on, we could arrive at a point where none of left/right options satisfy the condition $\in L \Rightarrow$ THIS means that before we made the wrong choice (i.e. both were fine, we chose the wrong \Rightarrow WE NEED TO BACKTRACKING).

BACKTRACKING \rightarrow IF we arrived at a dead point, we need to go back to previous attempts to modify them. It goes hand-in-hand with exhaustive search algorithms.

PLACE \rightarrow algorithm. \rightarrow segue proprio questo ragionamento.

If L empty \rightarrow return \Rightarrow abbiamo finito.

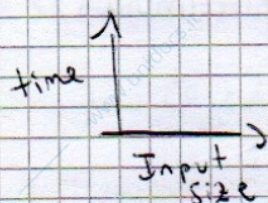
Remove y from X & add length $\Delta(y, X)$ to L

Costiamo togliendo y & lo raggio siamo a quella da analizzare.

Experimentally, this algorithm is much faster than the previous ones (at maximum 2 options each time) **Dot**:

WORST CASE SCENARIO: again exponential, BUT in practice it runs faster than the others \rightarrow

Data fitting → experimental analysis.



• Suppose we analytically know that the complexity is linear $O(n)$.

↳ How to find a linear function $f: \mathbb{N} \rightarrow \mathbb{N}$?

We can use a statistical approach: least-square method

INTUITION: We want to find a & b for the linear function.

We want to minimize the squared errors → Numpy & Scipy are libraries having it.

xdata = numpy.array([10.0, 20.0, 30.0, 40.0, 50.0])

ydata = numpy.array([7.3, 8.46, 10.2, 11.4, 13.08])

time associated to input size.

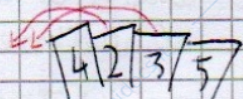
the estimate of our curfit → Ce la da l'algoritmo.

→ We will then plot it in LaTeX.

↳ we will use pgf

~~the plot of the data and the fit~~

• Why in InsertionSort there's return & e capo l'altro comando?
Let's assume we're given a set of cards.



Sorted AREA Unsorted AREA.

focusing on the first element of the unsorted area & moving it in the right position

$$t = a \cdot x^2 + b \cdot x$$

$$4.1279e^{-6} \cdot x^2 + (-1.2548e^{-3}) \cdot x$$

• since in same we have 20 & 25 coins \rightarrow the better to have 40 would be 2×20 & not $1 \times 25 + 1 \times 10 + 1 \times 5$

NOT CORRECT because the choices are optimal now BUT NOT always.

\hookrightarrow it's a greedy algorithm \rightarrow when you have a choice point you take the greedy (=the best now).

In this case it's better to start from the trivial solution \rightarrow we consider all the possible (no logic) combinations of coin that reach M & find the best solution.

\hookrightarrow BRUTE FORCE CHANGE \rightarrow you don't go completely blind (in line 2 there's some logic)

We give an UPPER BOUND

I focus on a small space: the combinations that do NOT go more than M

& since we go "blindly" there's NO change we miss any values

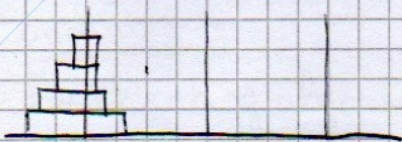
EXHAUSTIVE SEARCH

\hookrightarrow it is correct bc I spot the result I want & note it down.

08/03/2021

RECURSION = when the algorithm calls itself (partitioning the input into smaller parts which can be solved, then recombine the solutions).

HANOI PUZZLE



\rightarrow move all disks from peg 1 to peg 2, moving only one disk at the time without putting a bigger one on a smaller one.

Intuition: move the first $n-1$ disks on the 2nd peg and the last one on the 3rd. The solution with 3 disks is easy, if we do it w/ 4, we can reconduct ^{partially} the problem to the solution with 3 disks.

• So, what is the base CASE? Part executed without (recursive) recursive calls, it is correct when we only have one disk.

unused peg \leftarrow 6 - from peg - to peg (mathematical function to determine which of the 3 is actually not used)

• What is the recursive case? part which resorts the recursive calls.

ALGORITHMS & DATA STRUCTURES

04/03/2021

lectures will be held online → live a laudo, insegna al King's college

2 MODULES:

1. THEORY → 16 lectures of 2H

2. PRACTICALS → "hands-on" → we need to invent ways on how to solve problems

ABSTRACTION → understanding the structure of a problem from an abstract point of view, so that we can tackle the design.

We will to design how to solve problems → we will have a context of bioinformatics

Structural complexity → the complexity of a problem, NOT of an algorithm.

↳ to understand if

- our algorithm is not fast bc of the implementation
- the problem is intrinsically long to solve, no much I can do.

The point is that we don't want only "an answer", we want the best one.

EXAMS

① module: written (exercises) & oral (on how you reason)

② exercises during the practicals → 2/3 assignments → 10/14 days to solve it.

↳ scrivi soluzioni & then short report

• We will see what problems & algorithms are, the complexity of algorithms (this is important for deciding the algorithm to be used).

• How to design algorithms → following some design principles

— **Algorithms** are ways for solving problems.

— "Why should I learn Algorithmics?" → slide 7

↳ Usually we encounter problems for which people have already provided solutions (you simply import) → Sometimes, especially in difficult scenarios, is NOT like that → **You** need to deal w/ it. →

• There are \neq sets of points X , whose ΔX is the same.

↳ INTUITION: we could shift all the values by a n (ex. 2), & the pairwise distance would be kept.

↳ bc we're given 2 points, defining a distance but we're not told where these do start.

• With this algorithm we generate all the possible combinations $\binom{M-1}{n-2}$ so it's very slow.

= from a set of $M-1$ elements, need to take n

↓ We can improve a bit the performance

↓ Would it make sense to choose a point 1 in $\mathbb{1}^P$? The correctness would be impacted (BUT) the point is that we're wasting time trying conditions that **FOR SURE** are wrong (bc distance = 2 lost)

↳ So we can update the brute force algorithm, to produce another one. It is the same (BUT) with a difference in line 2, that the point must be contained in L .

$$\binom{M}{n-2}$$

combinations I generate

$$= O(n^{2n-4})$$

still exponential

$$O(M^{n-2})$$

for the previous

we choose this bc in principle $n \ll M$

here's example.

$$L = \{2, 998, 1000\}$$

$$\begin{cases} n = 3 \\ M = 1000 \end{cases}$$

• These algorithms are BOTH slow, even though the 2nd one is a bit less slow.

→ Before we started defining the extreme points (the initial ato) & $M = 10$. Then we said it would not make sense to put a point in L if so we'd observe a distance not present in the multiset input. last one can guide us in choosing the points.

what we can do: we ignore 10 (we already dealt w/ it) & we consider the 2nd biggest n ⇒ i.e. $\mathbb{8}$. Since the largest remaining distance is 8, there must be a point at distance 8 (0-8) but if



$L = \{10, 8, \dots\}$ → biggest element → To justify: there are 2 points of 10.

In Python arrays are expressed as lists.

↓
here elements accessed as elements in list. → the index = 0 of the first element.

In pseudocode → just indexing → the index of first element = 1

Problem → Change Problem

quarters → 25 cents

dimes → 10 cents

nicks → 5 cents

pennies → 1 cent

INPUT: M → amount of money in cents.

OUTPUT: ↓
we want to give an amount of cents, whose sum is M but minimized in number.

- ① We start giving as many quarters as we can ($\leq M$).
Then we pass to a less valuable coin.

HOW TO CHECK CORRECTNESS?

- ① we can test the algorithm on a huge set of inputs. The problem is that we're sure it works on those → se likely correct, BUT there may be an input, not tested, which gives an error

- ② we can prove the correctness.

↳ it is correct if on all possible inputs, the algorithm always gives the right answer.

↳ when is it not correct? the negation! & the negation is that at least ① input does not give the right answer.

generalized.

BetterChange(M, c, d):

array of values in decreasing order

1 $r \leftarrow M$

remainder → i soldi che mi mancano da dare.

2 for $k \leftarrow 1$ to d

3 $i_k \leftarrow r / c_k$ → integer division ex. $\frac{67}{50} \rightarrow 1$

4 $r \leftarrow r - c_k \cdot i_k$ → $r = 17$ → this enters again cannot be divided by coin of value 20, but yes with 10.

5 return (i_1, i_2, \dots, i_d)

||
NOT CORRECT bc in some situations it does not work

- Practical point of view

EXPERIMENTAL APPROACH

slides

→ How can we do it? By profiling a program → in slide spiegato.

To use it we will use cProfile (you simply import it).

In the example of Fibonacci we have:

- the import statement
- cProfile.run('fibonacci(12)')

through this we're telling the system to perform Fibonacci, passing the value '12' & keep track of the timing & statistics.

→ We have 468 calls & the time 0.00, bc it was very fast.

we call the Statistics for Fibonacci basically.

In the next example:
48315633 → Fibonacci calls
48315636 →

tells you how many times each function was called

We're interested in Cumtime → computation time: the time spent to execute from the first call.

ASSIGNMENTS (some will be marked)

① Describe the situation (streamline & formalize) → loro te lo danno scritto in generico → leggi la descrizione, come posso renderlo → what's the input/output / how I proceed.

② You will have to implement in Python → loro possono chiederti anche la complessità → you draw a graph & 2 lines of conclusions.

explain why correct & what you think of the complexity / draw a graph, spiega. → PAROLE IN CONCLUSIONE

→ DA USARE LATEX. → tailored for scientific research.
Lo X scrivere il report

www.overleaf.com

Crea profilo → FREE ACCOUNT

→

• BRUTE FORCE MOTIF SEARCH AGAIN \rightarrow basically here we integrated the previous algorithms, but we're looking at all leaves of the tree. Still we haven't cut parts of it, \rightarrow we need to figure out how to search in the tree. to explore

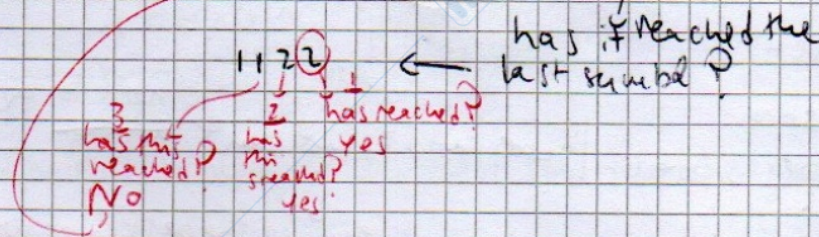
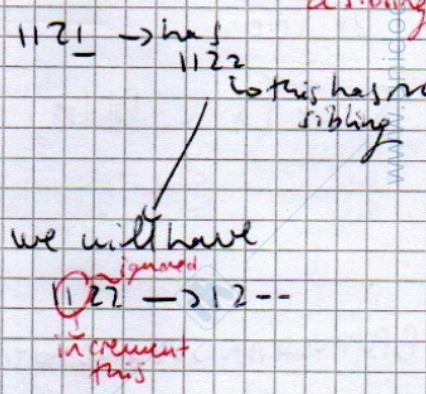
\rightarrow DEPTH FIRST SEARCH: we start exploring in depth \rightarrow we go from the root (1) to the 5 \rightarrow the leaf \rightarrow then the sibling; Now we go back on our way to 7, then again in depth (8 & 9) then up to 10, then in depth. First we explore a node, then left ^{end} side sub-tree and then the right ^{end} side sub-tree.

\checkmark we go to the maximum depth as we can; once we arrive to a leaf (5). The following node analyzed is the sibling; after this there are no more siblings. Once we ~~don't~~ have arrive at 9, the next node to analyze is the first sibling NOT yet explored of an ancestor of the current case. 9 \rightarrow 10 case.

16 \rightarrow 17
 \hookrightarrow all the left end side has been explored

NEXT VERTEX (a, i, L, k)
 \uparrow current node level at the moment

- 1 if $i < L$ \rightarrow if so, we're not on a leaf, we want to go down. (12--)
- 2 $a_{i+1} \leftarrow 1$ (12+)
- 3 return (a, i+1)
- 4 else \rightarrow we're in a leaf.
- 5 for $j \leftarrow L$ to i ^{to last position, to first last one.}
- 6 if $a_j < k$ ~~stop~~ ^{we look at the last position & ask: is this the last symbol? If not, it means there's a sibling}
- 7 $a_j \leftarrow a_j + 1$
- 8 return (a, 0) ^{once we've finished, we return 0 as level of exploration.}
- 9 return (a, j) ^{our level}



• SIMPLER MOTIF SEARCH

line \rightarrow 5/6 \rightarrow manage the exploration when we're at mid-level \rightarrow we simply generate the next vertex.

\hookrightarrow so basically we're not doing anything

ENUMERATING SETS in Python

31/03/2021

- Before to implement the algorithm \rightarrow try drawing ideas on paper.
- IDEA \rightarrow we have to build a subset of a given set \rightarrow you start from an empty set \rightarrow you encounter the 1st element \rightarrow to add or not? & so on \rightarrow & you draw all possible combinations.

LAB 5 We're going to learn a trick for not running out of memory! \rightarrow GENERATORS

16/04/21

How to deal with sets in Python

- can be presented with curly braces $\{ \dots \}$
- operator `in`
- the difference w/ lists is that even if we have more occurrences of the same element in the representation (as in the mathematical concept), the element is printed only once.

! we cannot use $\{ \}$ bc the means empty dictionary

`set()` \rightarrow for the empty set

UNION:

$\{0, 1, 2, 3, 4\} \cup \{1, 3, 5, 7\}$

INTERSECTION:

&

DIFFERENCE:

$\{0, 1, 2, 3\} - \{1, 3, 5, 7\}$

$\{0, 2\}$

\hookrightarrow basically we remove from the first set the elements in the second.

non modifying operators

\hookrightarrow the operands are not modified, the operator generates a new

`.add()`

`.remove()`

`|=` = live plus equal

\rightarrow Sometimes it is useful to convert sets into lists & viceversa:
When don't know the order of elements for sets

• We have `.copy()`

• we have the command `len()`

GENERATORS

\rightarrow In the ex. we did, when we ~~all set~~ generated all the possible combinations; BUT THE CONCEPT IS THAT

• APPROXIMATION RATIOS \rightarrow how much our greedy algorithm is far from the optimal solution.

\rightarrow the bigger the n size, the more the approx. ratio \Rightarrow for our algorithm, so approx. ratio not fixed in our case.

Today we'll see an algorithm whose ALGORITHM RATIO is constant. (it will be 4), that is not as an approx. ratio of error.

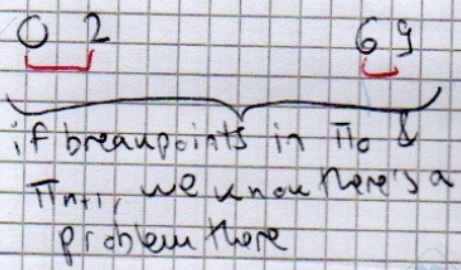
what we do! GIVING NAMES

• when we consider a permutation $\pi = \pi_1, \dots, \pi_n$, we introduce 2 more numbers

$\pi_0 = 0$ & $\pi_{n+1} = n+1$
tail

• adjacency \rightarrow one element is the successive of the other \rightarrow we don't care if 1-2 or 2-1 (bigger first of later) $\Rightarrow |\pi_i - \pi_{i+1}| = 1$

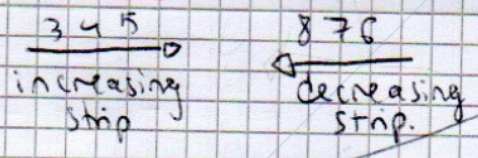
• breakpoint \rightarrow 1-3 or 5-8, when they're not consecutive
 \hookrightarrow essentially intuitively tell us about problems



INTUITIVELY: the more breakpoints we have the faster from the sol. we are

the only thing is that π_0 & π_{n+1} ARE NOT counted for the reversal

• a Strip = a sequence of adjacencies \rightarrow any interval between 2 breakpoints



If many ~~breakpoints~~ we need to have many reversals, in the end there will be no breakpoints.

MAX n° : $n+1$ if completely unsorted seq.

n° of breakpoints: measure of how far we are.

What does it happen to a breakpoint when we reverse smth?

• Assume we're going to reverse 8 7 6, what happens is that everything in the interval is reversed \rightarrow everything that was an adjacency remains an adjacency. We're just changing the order of things BUT NOT the order of elements in the segment. Whatever it was a breakpoint, will remain a breakpoint as well.

When we do a reversal, the n° of adjacencies & breakpoints remains the same WITHIN the segment.

! What changes is what happens at the borders \rightarrow in the best scenario the n° of breakpoints could diminish by 2. \rightarrow (! at most \rightarrow in the best scenario)