

Predicting Solvent Accessibility of Protein Residues using High Order Conditional Random Fields

Dongbo Bu

Institute of Computing Technology
Chinese Academy of Sciences, Beijing, China

The origin of the word “Algorithm”

Figure: Muhammad ibn Musa al-Khwarizmi (C. 780—850), a Persian scholar, formerly Latinized as Algoritmi

- In the twelfth century, Latin translations of his work on the Indian-Arabic numerals introduced the decimal positional number system to the Western world.

- Al-Khwarizmi's *The Compendious Book on Calculation by Completion and Balancing* presented the first systematic solution of **linear** and **quadratic equations** in Arabic.
- Two words:
 - **Algebra**: from Arabic "al-jabr" meaning "reunion of broken parts" — one of the two operations he used to solve equations
 - **Algorithm**: a step-by-step set of operations to get solution to a problem

Algorithm design: the art of computer programming

Our philosophy on the design and exposition of algorithms is nicely illustrated by the following analogy with an aspect of Michelangelos's art:

*A major part of his effort involved looking for interesting pieces of stone in the quarry and staring at them for long hours **to determine the form they naturally wanted to take**. The chisel work exposed, in a minimal manner, this form.*

By analogy, we would like to start with a clean, simply stated problem.

*Most of the algorithm design effort actually goes into **understanding the algorithmically relevant combinatorial structure of the problem.***

*The algorithm exploits this structure in a minimal manner.....
with emphasis on stating the structure offered by the problems,
and keeping the algorithms minimal.*

(See extra slides.)

Basic algorithmic techniques

- **Divide-and-conquer**: Let's start from the “smallest” problem first, and investigate whether a large problem can **reduce to smaller subproblems**.
- **Improvement**: Let's start from **an initial complete solution**, and try to improve it step by step.
- **“Intelligent” enumeration**: we enumerate **all possible complete solutions**, but employ some techniques to prune the search tree.

The first example: calculating the greatest common divisor (gcd)

The first problem: calculating gcd

Definition (gcd)

The greatest common divisor of two integers a and b , when at least one of them is not zero, is the largest positive integer that divides the numbers without a remainder.

- Example:
 - The divisors of 54 are: 1, 2, 3, 6, 9, 18, 27, 54
 - The divisors of 24 are: 1, 2, 3, 4, 6, 8, 12, 24
 - Thus, $\text{gcd}(54, 24) = 6$.

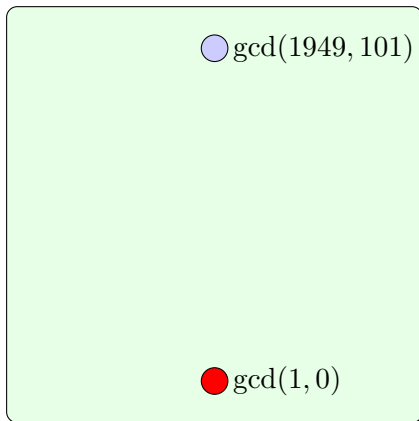
The first problem: calculating gcd

INPUT: two n bits numbers a , and b ($a \geq b$)

OUTPUT: $\gcd(a, b)$

- Observation: the problem size can be measured by using n ;
- Let's start from the “smallest” instance: $\gcd(1, 0) = 1$;
- But how to efficiently solve a “larger” instance, say $\gcd(1949, 101)$?

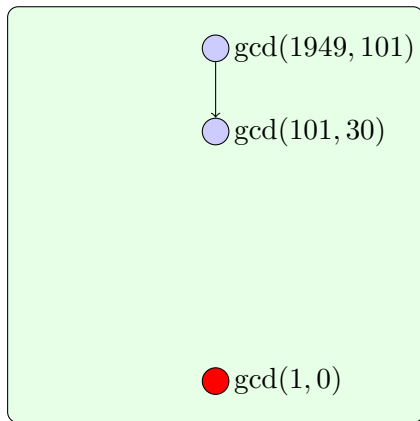
Problem instances



- Observation: a large problem can reduce to a smaller subproblems:
- $\text{gcd}(1949, 101) = \text{gcd}(101, 1949 \bmod 101) = \text{gcd}(101, 30)$

Strategy: reduce to “smaller” problems

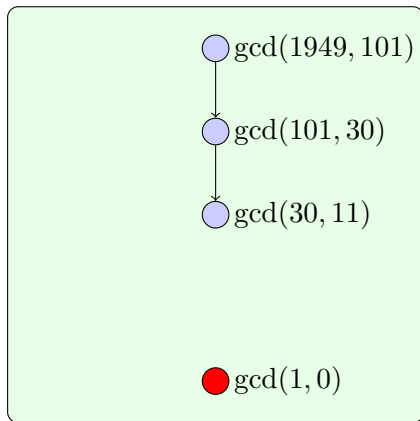
Problem instances



- $\text{gcd}(101, 30) = \text{gcd}(30, 101 \bmod 30) = \text{gcd}(30, 11)$

Strategy: reduce to "smaller" problems

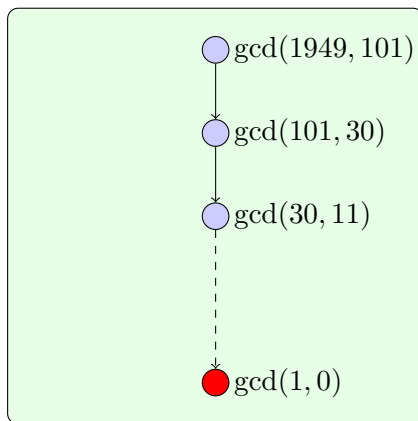
Problem instances



- $\text{gcd}(30, 11) = \text{gcd}(11, 30 \bmod 11) = \text{gcd}(11, 8)$

Strategy: reduce to "smaller" problems

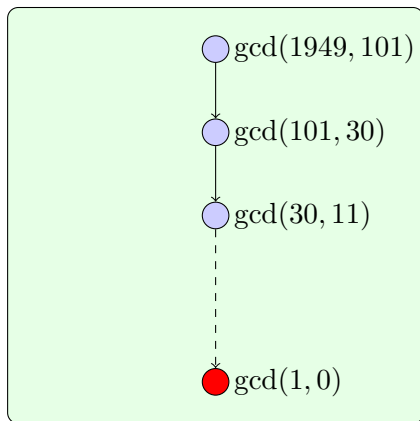
Problem instances



- $\gcd(30, 11) = \gcd(11, 8) = \gcd(8, 3) = \gcd(3, 2) = \gcd(2, 1) = \gcd(1, 0) = 1$

Sub-instance relationship graph

Problem instances



- Node: subproblems
- Edge: reduction relationship

Euclid algorithm

```
1: function Euclid( $a, b$ )  
2: if  $b = 0$  then  
3:   return  $a$ ;  
4: end if  
5: return Euclid( $b, a \bmod b$ ) ;
```


Theorem

Suppose a is a n -digit integer. $\text{Euclid}(a, b)$ ends in $O(n^3)$ time.

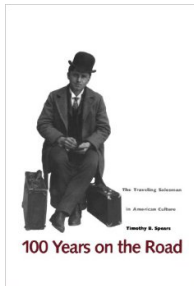
Proof.

- There are at most $2n$ recursive calling.
 - Note that $a \bmod b < \frac{a}{2}$.
 - After two rounds of recursive calling, both a and b shrink at least a half size.
- At each recursive calling, the \bmod operation costs $O(n^2)$ time.



The second example: travelling salesman problem (TSP)

TSP: a concrete example



- In 1925, H. M. Cleveland, a salesman of the Page seed company, traveled 350 cities to gather order form.
- Of course, the shorter the total distance, the better.

How did they do? Pin and wire!

OFFICE APPLIANCES

201

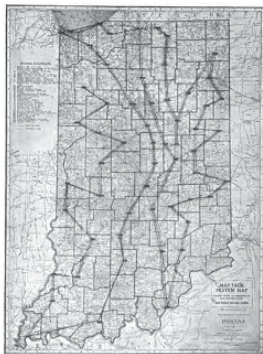


MAP CABINET

Courtesy of Rand-McNally

202

SECRETARIAL STUDIES



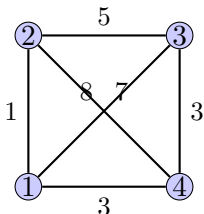
MAP SHOWING ROUTING OF SALESMEN BY PINS AND CORDS
Courtesy of Rand-McNally

- Two pictures excerpted from *Secretarial Studies*, 1922.

TRAVELLING SALESMAN PROBLEM

INPUT: a list of n cities $V = \{1, 2, \dots, n\}$, and a distance matrix D , where d_{ij} ($1 \leq i, j \leq n$) denotes the distance between city i and j

OUTPUT: the shortest tour that visits each city exactly once and returns to the origin city

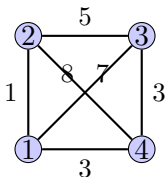


#Tours: 6

- Tour 1: $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 1$ (distance: 12)
- Tour 2: $1 \rightarrow 2 \rightarrow 4 \rightarrow 3 \rightarrow 1$ (distance: 21)
- Tour 3: $1 \rightarrow 3 \rightarrow 2 \rightarrow 4 \rightarrow 1$ (distance: 23)

Trial 1: divide and conquer

Consider a tightly related problem



Definition

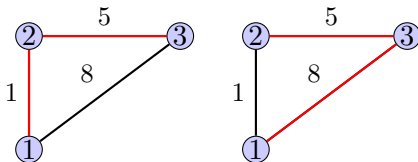
$D(S, e)$ = the minimum distance, starting from city 1, visiting all cities in S once and exactly once, and finishing at city e .

- There are 3 cases of the city from which we return to 1.
- Thus, the shortest tour can be calculated as:

$$\begin{aligned} & \min\{d_{2,1} + D(\{3, 4\}, 2), \\ & d_{3,1} + D(\{2, 4\}, 3), \\ & d_{4,1} + D(\{2, 3\}, 4)\} \end{aligned}$$

Consider the smallest problem

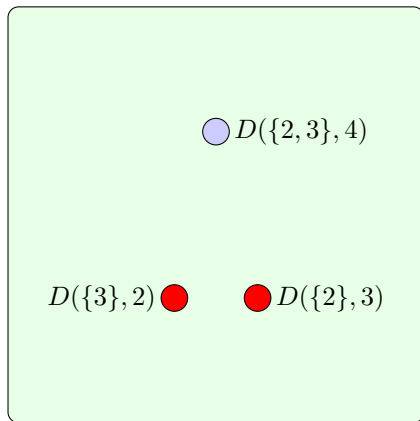
- It is trivial to calculate $D(S, e)$ when S consists of only 1 city.



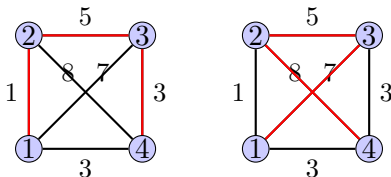
- $D(\{2\}, 3) = d_{12} + d_{23}$ $D(\{3\}, 2) = d_{13} + d_{32}$
- But how to solve a larger problem, say $D(\{2, 3\}, 4)$?

Sub-instance relationship graph

Problem instances



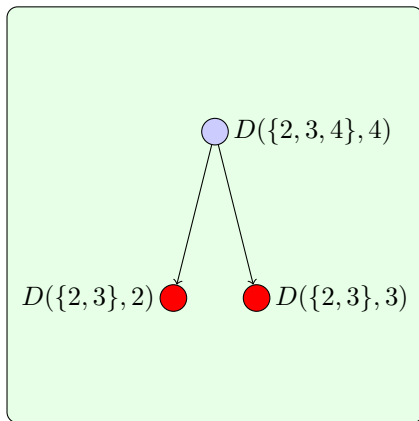
Divide a larger problem into smaller problems



- $D(\{2, 3\}, 4) = \min\{d_{34} + D(\{2\}, 3), d_{24} + D(\{3\}, 2)\}$

Reduce to smaller problems

Problem instances



Held-Karp algorithm [1962]

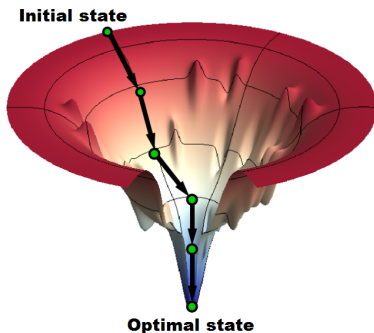
```
1: function  $TSP(V, D)$   
2: return  $\min_{e \in V, e \neq 1} D(V - \{e\}, e) + d_{e1};$ 
```

```
1: function  $D(S, e)$   
2: if  $S = \{v\}$  then  
3:    $D(S, e) = d_{1v} + d_{ve};$   
4:   return  $D(S, e);$   
5: end if  
6: return  $\min_{i \in S, i \neq e} D(S - \{i\}, i) + d_{ei};$ 
```

- Time complexity: $O(2^n n^2)$.

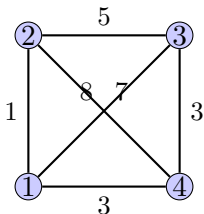
Trial 2: Improvement strategy

Solution space



- Node: a complete solution. Each node is associated with an objective function value.
- Edge: if two nodes are neighbors, an edge is added to connect them. Here "neighbours" refers to two nodes with small difference.
- Improvement strategy: start from an initial solution, and try to improve it step by step.

IMPROVEMENT strategy



- Note that a **complete solution** can be expressed as a permutations of the n cities;
- Let's start from **an initial complete solution**, and try to improve it;

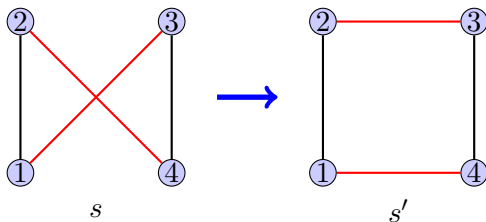
IMPROVEMENT strategy

```
1: Let  $s$  be an initial tour;
2: while TRUE do
3:   select a new tour  $s'$  from the neighbourhood of  $s$ ;
4:   if  $s'$  is shorter than  $s$  then
5:      $s = s'$ ;
6:   end if
7:   if stopping( $s$ ) then
8:     return  $s$ ;
9:   end if
10: end while
```

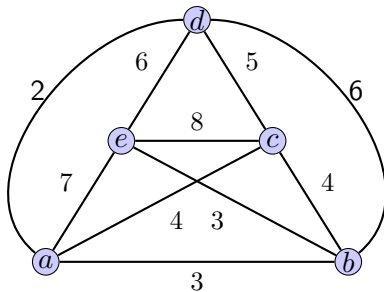
Here, **neighbourhood** is introduced to describe how to change an existing tour into a new one;

But how to define neighbourhood of a tour?

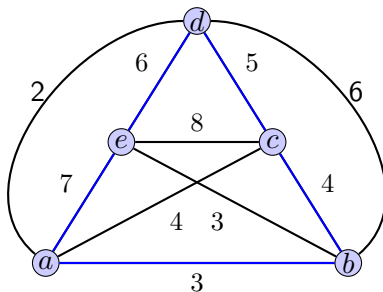
- 2-opt strategy: if s' and s differ at only two edges (Note: 1-opt is impossible)



An example

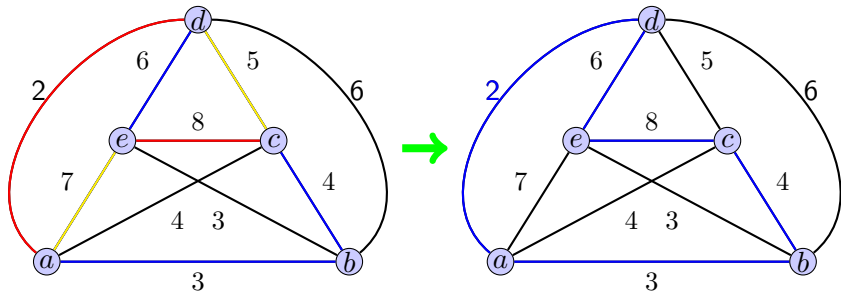


Initial tour s



Initial complete solution s : $a \rightarrow b \rightarrow c \rightarrow d \rightarrow e \rightarrow a$ (distance: 25)

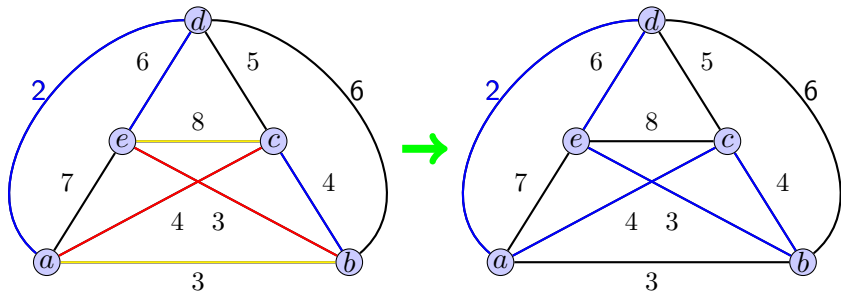
A 2-opt operation: $s \Rightarrow s'$



Initial solution s : $a \rightarrow b \rightarrow c \rightarrow d \rightarrow e \rightarrow a$ (distance: 25)

Improve from s to s' : $a \rightarrow b \rightarrow c \rightarrow e \rightarrow d \rightarrow a$ (distance: 23)

One more 2-opt operation: $s' \Rightarrow s''$

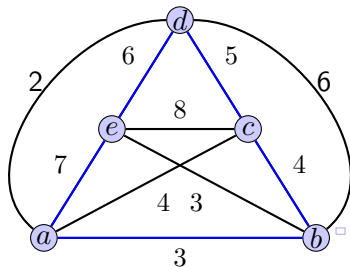


- A complete solution s' : $a \rightarrow b \rightarrow c \rightarrow e \rightarrow d \rightarrow a$ (distance: 23)
- Improve from s' to s'' : $a \rightarrow c \rightarrow b \rightarrow e \rightarrow d \rightarrow a$ (distance: 19)
- Done! No 2-OPT can be found to improve further.

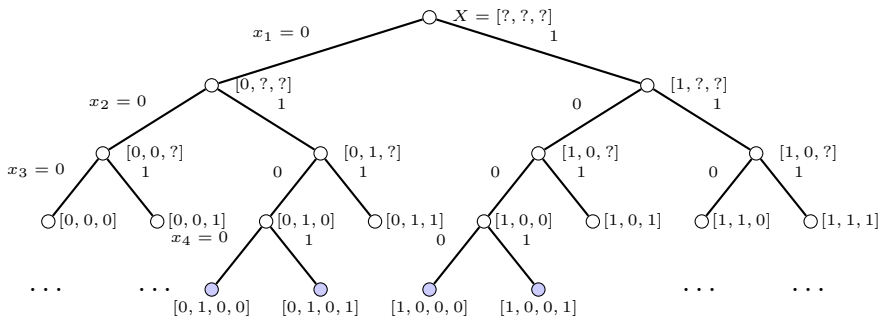
Trial 3: “Intelligent” enumeration strategy

Solution form

- Note that a complete solution can be expressed as a sequence of $n - 1$ edges. Given a certain order of the m edges, a complete solution can be represented as $X = [x_1, x_2, \dots, x_m]$, $x_i = 0/1$. If the edge i was used, we set $x_i = 1$, and otherwise $x_i = 0$.
- For example, the tour $a \rightarrow b \rightarrow c \rightarrow d \rightarrow e \rightarrow a$ can be represented as $X = [1, 0, 0, 1, 1, 0, 0, 1, 0, 1]$ under the order $e_1 = \langle a, b \rangle, e_2 = \langle a, c \rangle, e_3 = \langle a, d \rangle, e_4 = \langle b, c \rangle, e_5 = \langle b, d \rangle, e_6 = \langle b, e \rangle, e_7 = \langle c, d \rangle, e_8 = \langle c, e \rangle, e_9 = \langle d, e \rangle$.



Partial solution tree

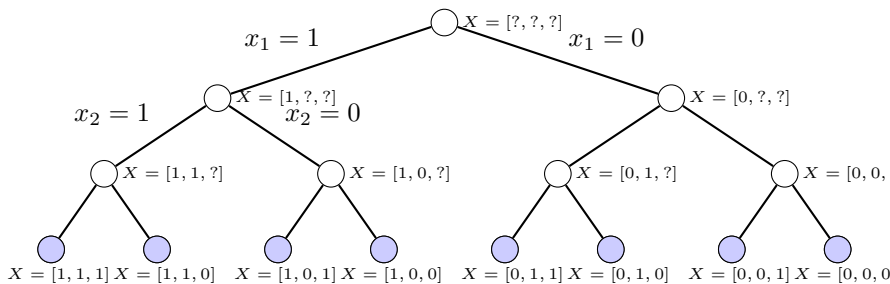


- Internal node: partial solution
- Leaf node: complete solution

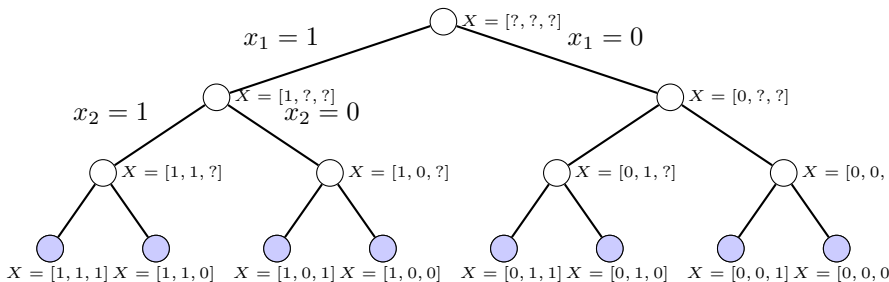
Enumerate all complete solutions: Backtrack algorithm

```
1: Start with the original problem  $P_0$ ;  
2: Let  $A = \{P_0\}$ . Here  $A$  denotes the active subproblems.  
3:  $bestsofar = \infty$ ;  
4: while  $A \neq NULL$  do  
5:   choose a subproblem  $P \in A$ , and remove it from  $A$ ;  
6:   expand  $P$  into smaller subproblems  $P_1, P_2, \dots, P_k$ ;  
7:   for  $i = 1$  to  $k$  do  
8:     if  $P_i$  corresponds to a complete solution then  
9:       update  $bestsofar$ ;  
10:    else  
11:      insert  $P_i$  into  $A$ ;  
12:    end if  
13:  end for  
14: end while  
15: return  $bestsofar$ 
```

An example

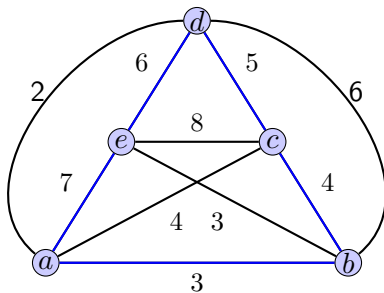


Question: can we make the enumeration smart?



- Basic idea: only complete solutions are associated with objective function value. Could we estimate objective function value for partial solutions?
- Alternative viewpoints of **partial solution**
 - A partial solution represents a sub-problem.
 - A partial solution represents a set of complete solutions.
- Let's use the **lower bound** of these complete solutions as an estimation of the quality of a partial solution.

Estimate quality of partial solution: an example



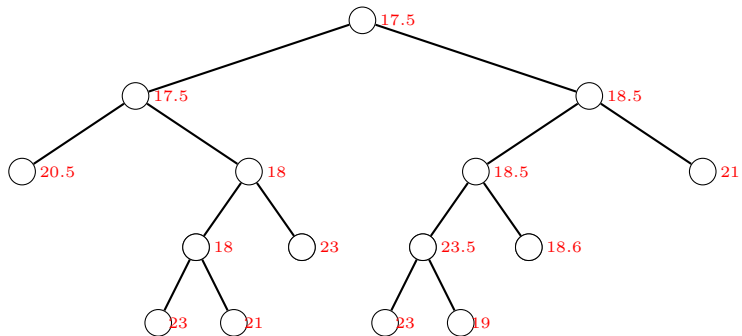
For the partial solution $X = [?, ?, ?]$, we estimate the shortest tour as below:

- for each node, we select the shortest two adjacent edges
- the sum of these $2n$ edges is less than 2 times the optimal tour
- thus, we can give a lower bound as $\frac{1}{2}(5 + 6 + 7 + 8 + 9) = 17.5$

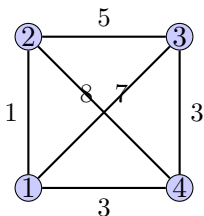
“Intelligent” enumeration

```
1: Start with the original problem  $P_0$ ;  
2: Let  $A = \{P_0\}$ . Here  $A$  denotes the active subproblems.  
3:  $bestsofar = \infty$ ;  
4: while  $A \neq NULL$  do  
5:   choose a subproblem  $P \in A$ , and remove it from  $A$ ;  
6:   expand  $P$  into smaller subproblems  $P_1, P_2, \dots, P_k$ ;  
7:   for  $i = 1$  to  $k$  do  
8:     if  $P_i$  corresponds to a complete solution then  
9:       update  $bestsofar$ ;  
10:    else  
11:      if  $lowerbound(P_i) \leq bestsofar$  then  
12:        insert  $P_i$  into  $A$ ;  
13:      end if  
14:    end if  
15:  end for  
16: end while  
17: return  $bestsofar$ 
```

An example

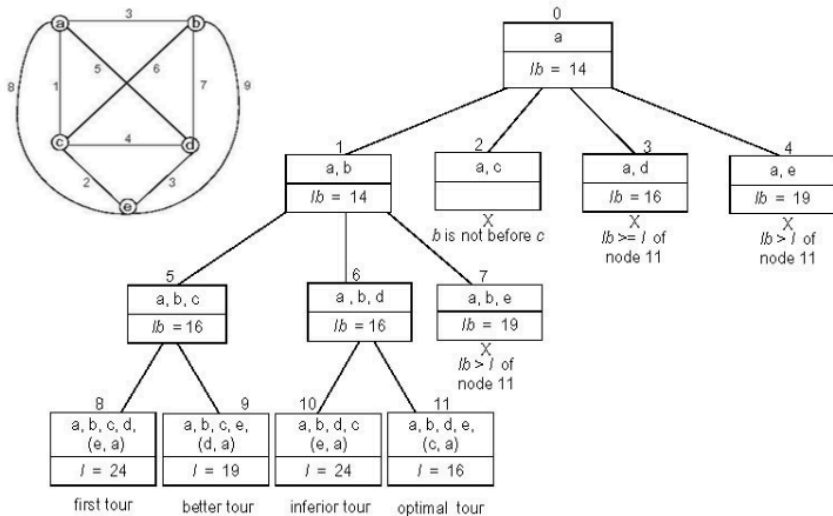


BACKTRACKING strategy: another trial



- Note that a complete solution can be expressed as a sequence of n nodes;
- Formally, we can represent a complete solution as:
$$X = [v_1, v_2, \dots, v_n].$$

Enumeration tree and backtracking



Backtracking

- 1: Start with the original problem P_0 ;
- 2: Let $A = \{P_0\}$. Here A denotes the active subproblems.
- 3: $bestsofar = \infty$;
- 4: **while** $A \neq NULL$ **do**
- 5: **choose** a subproblem $P \in A$, and remove it from A ;
- 6: **expand** P into smaller subproblems P_1, P_2, \dots, P_k ;
- 7: **for** $i = 1$ to k **do**
- 8: **if** (P_i) corresponds to a complete solution **then**
- 9: update $bestsofar$;
- 10: **end if**
- 11: **if** $lowerbound(P_i) \leq bestsofar$ **then**
- 12: insert P_i into A ;
- 13: **end if**
- 14: **end for**
- 15: **end while**

Time complexity and space complexity

Time complexity and space complexity

- Time (space) complexity of an algorithm quantifies the time (space) taken by the algorithm.
- Since the time costed by an algorithm grows with the size of the input, it is traditional to describe running time as a function of the input size.
 - **input size**: The best notation of input size depends on the problem being studied.
 - For the TSP problem, the **number of cities in the input**.
 - For the MULTIPLICATION problem, the **total number of bits** needed to represent the input number is the best measure.

Running time: we are interested in its growth rate

- A straightforward way is to use the exact seconds that a program used. However, this measure highly depends on CPU, OS, compiler, etc.
- Several simplifications to ease analysis of Held-Karp algorithm:
 - ① We simply use the number of primitive operations (rather than the exact seconds used) under the assumption that a primitive operation costs constant time. Thus the running time is $T(n) = an^2 + bn + c$ for some constants a, b, c .
 - ② We consider only the leading term, i.e. an^2 , since the lower order terms are relatively insignificant for large n .
 - ③ We also ignore the leading term's coefficient a since it is less significant than the growth rate.
- Thus, we have $T(n) = an^2 + bn + c = O(n^2)$. Here, the letter O denotes **order**.

Big O notation

- Recall that big O notation is used to describe the **error term** in Taylor series, say:

$$e^x = 1 + x + \frac{x^2}{2} + O(x^3) \text{ as } x \rightarrow 0$$

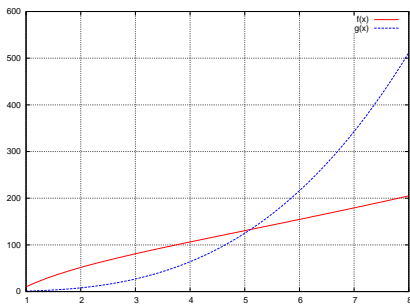


Figure: Example: $f(x) = O(g(x))$ as there exists $c > 0$ (e.g. $c = 1$) and $x_0 = 5$ such that $f(x) < cg(x)$ whenever $x > x_0$

Big Ω and Big Θ notations

- In 1976 D.E. Knuth published a paper to justify his use of the Ω -symbol to describe a stronger property. Knuth wrote: "For all the applications I have seen so far in computer science, a stronger requirement [...] is much more appropriate".
- He defined

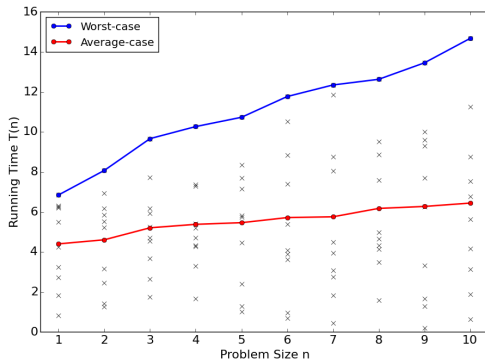
$$f(x) = \Omega(g(x)) \Leftrightarrow g(x) = O(f(x))$$

with the comment: "Although I have changed Hardy and Littlewood's definition of Ω , I feel justified in doing so because their definition is by no means in wide use, and because there are other ways to say what they want to say in the comparatively rare cases when their definition applies".

- Big Θ notation is used to describe " $f(n)$ grows asymptotically as fast as $g(n)$ ".

$$f(x) = \Theta(g(x)) \Leftrightarrow g(x) = O(f(x)) \text{ and } f(x) = O(g(x)).$$

Worst case and average case



- Worst-case: the case that takes the longest time;
- Average-case: we need know the distribution of the instances;