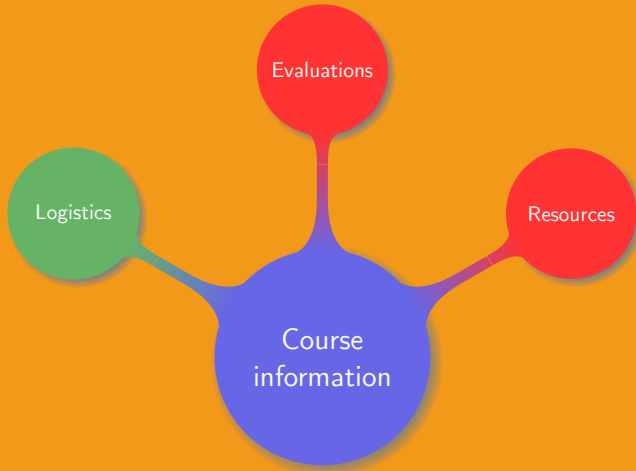


Introduction to Algorithms

Manuel – Fall 2019

0. Course information



Teaching team:

- Instructor: Manuel (charlem@sjtu.edu.cn)
- Teaching assistants:
 - Yifei (zhangyifei-chelsea@sjtu.edu.cn)
 - Niyiqiu (lnyq10@sjtu.edu.cn)

Teaching team:

- Instructor: Manuel (charlem@sjtu.edu.cn)
- Teaching assistants:
 - Yifei (zhangyifei-chelsea@sjtu.edu.cn)
 - Niyiqiu (lnyq10@sjtu.edu.cn)

Important rules:

- When contacting a TA for an important matter, CC the instructor
- Prepend [VE477] to the subject, e.g. Subject: [VE477] Grades
- Use [SJTU jBox service](#) to share large files (> 2 MB)

Teaching team:

- Instructor: Manuel (charlem@sjtu.edu.cn)
- Teaching assistants:
 - Yifei (zhangyifei-chelsea@sjtu.edu.cn)
 - Niyiqiu (lnyq10@sjtu.edu.cn)

Important rules:

- When contacting a TA for an important matter, CC the instructor
- Prepend [VE477] to the subject, e.g. Subject: [VE477] Grades
- Use [SJTU jBox service](#) to share large files (> 2 MB)

Never send large files by email

Course arrangements:

- Lectures:
 - Tuesday 8:00 – 9:40
 - Thursday 8:00 – 9:40
 - Friday 12:10 – 13:50 (even weeks)
- Office hours:
 - Tuesday 9:40 – 11:20
 - Thursday 9:40 – 11:20

Appointments outside of the office hours can be taken by email

Main goals of this course:

- Become familiar with the most common problems and algorithm paradigms
- Understand how to properly analyse and abstract a problem
- Be able to identify or design clear and efficient algorithms to solve a problem

Main goals of this course:

- Become familiar with the most common problems and algorithm paradigms
- Understand how to properly analyse and abstract a problem
- Be able to identify or design clear and efficient algorithms to solve a problem

Solve a given problem, assess the validity, quality, and efficiency of the solution

Learning strategy:

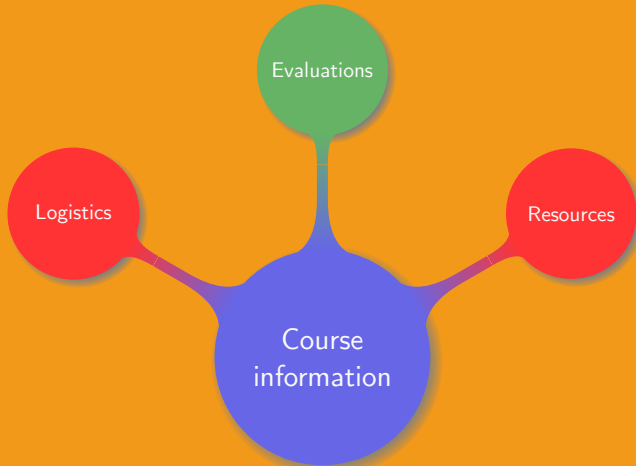
- Course side:
 - 1 Understand the basic concept of algorithmic
 - 2 Know the most common problems and their solutions
 - 3 Get an overview of the wide applications of algorithms

Learning strategy:

- Course side:
 - ① Understand the basic concept of algorithmic
 - ② Know the most common problems and their solutions
 - ③ Get an overview of the wide applications of algorithms
- Personal side:
 - ① Read and write code
 - ② Relate known strategies to new problems
 - ③ Perform extra research

Detailed goals:

- Be able to write clean and clear pseudocode
- Be proficient at using all the basic algorithm paradigms
- Be able to assess the difficulty of a given problem, e.g belongs to P, NP, PSPACE...
- Develop critical thinking abilities
- Know when and how to apply dynamic programming
- Know when and how to apply linear programming
- Have a precise idea of the pros and cons for the most common data structures
- Know how to efficiently solve the most common mathematical problems
- Have a basic idea on how to design multi-threaded algorithms
- Be able to efficiently implement the most common algorithms



Homework:

- Total: 8
- Content: basic concepts, critical thinking, prove results

Labs:

- Total: 8
- Content: implement common algorithms, learn Python

Project:

- Total: 1, split into three phases
- Content: write a catalog of the most common algorithms

Challenge:

Grade weighting:

- Assignments: 20%
- Projects: 20%
- Labs: 10%
- Midterm exam: 25%
- Final exam: 25%

Grade weighting:

- Assignments: 20%
- Projects: 20%
- Labs: 10%
- Midterm exam: 25%
- Final exam: 25%

Assignment submissions:

- Bonus: +10% for a work fully written in \LaTeX , bounded to 100%
- Penalty: -10% for a work not written in a neat and legible fashion
- Late policy: -10% per day, not accepted after three days

Grade weighting:

- Assignments: 20%
- Projects: 20%
- Labs: 10%
- Midterm exam: 25%
- Final exam: 25%

Assignment submissions:

- Bonus: +10% for a work fully written in \LaTeX , bounded to 100%
- Penalty: -10% for a work not written in a neat and legible fashion
- Late policy: -10% per day, not accepted after three days

Grades will be curved with the median in the range $\llbracket B, B+ \rrbracket$

General rules:

- Not allowed:
 - Reuse the code/work from other students
 - Reuse the code/work from the internet
 - Give too many details on how to solve an exercise

General rules:

- Not allowed:
 - Reuse the code/work from other students
 - Reuse the code/work from the internet
 - Give too many details on how to solve an exercise
- Allowed:
 - Share ideas and understandings on the course
 - Provide general directions on where or how to find information

Documents allowed during the exams:

- The lecture slides with **notes on them** (paper or electronic)
- A mono or bilingual dictionary

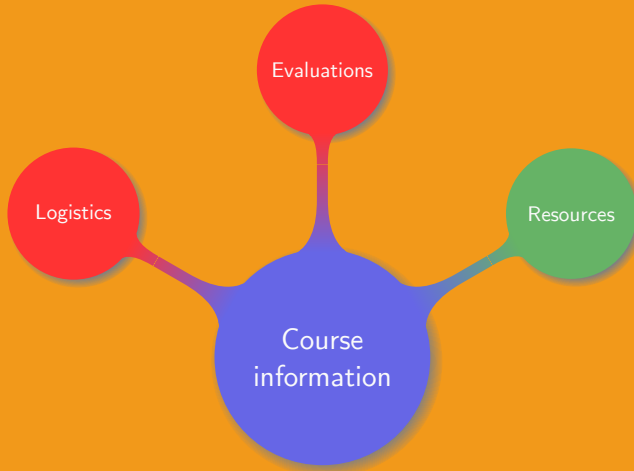
Group works:

- Every student in a group is responsible for his group submission
- If a student breaks the Honor Code, the whole group is sent to Honour Council

Contact us as early as possible when:

- Facing special circumstances (e.g. full time work, illness...)
- Feeling late in the course
- Feeling to work hard without any result

Any late request will be rejected



On **Canvas** platform:

- Course materials:
 - Syllabus
 - Lecture slides
 - Homework
 - Labs
 - Projects
 - Challenges
- Course information:
 - Announcements
 - Grades
 - Notifications
 - Polls

Places to find information:

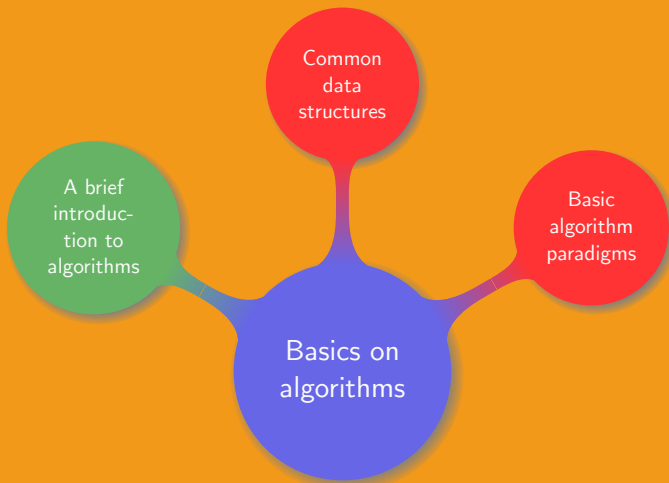
- *Algorithm Design*, J. Kleinberg and E. Tardos
- *Introduction to Algorithms*, H. Cormen, C. Leiserson, R. Rivest, and C. Stein
- *The Art of Computer Programming*, D. Knuth
- Piazza
- Search information online, i.e. $\{\text{websites} \setminus \{\text{local Chinese network}\}\}$

Places to find information:

- *Algorithm Design*, J. Kleinberg and E. Tardos
- *Introduction to Algorithms*, H. Cormen, C. Leiserson, R. Rivest, and C. Stein
- *The Art of Computer Programming*, D. Knuth
- Piazza
- Search information online, i.e. $\{\text{websites} \setminus \{\text{local Chinese network}\}\}$

Never use Baidu in any course

1. Basics on algorithms



An Algorithm is a recipe telling the computer how to solve a problem

An Algorithm is a recipe telling the computer how to solve a problem

Example.

Detail an algorithm to prepare a jam sandwich

Actions: cut, listen, spread, sleep, read, take, eat, dip

Things: knife, guitar, bread, honey, jam jar, sword

An Algorithm is a recipe telling the computer how to solve a problem

Example.

Detail an algorithm to prepare a jam sandwich

Actions: cut, listen, spread, sleep, read, take, eat, dip

Things: knife, guitar, bread, honey, jam jar, sword

Algorithm. (*Sandwich making*)

Input : 1 bread, 1 jamjar, 1 knife

Output: 1 jam sandwich

- 1 take the knife and cut 2 slices of bread;
 - 2 dip the knife into the jamjar;
 - 3 spread the jam on the bread, **using the knife**;
 - 4 assemble the 2 slices together, **jam on the inside**;
-

An algorithm systematically solves a *well-defined* problem:

- The *Input* is clearly expressed
- The *Output* solves the problem
- The *Algorithm* provides a precise step-by-step procedure starting from the Input and leading to the Output

An algorithm systematically solves a *well-defined* problem:

- The *Input* is clearly expressed
- The *Output* solves the problem
- The *Algorithm* provides a precise step-by-step procedure starting from the Input and leading to the Output

Algorithms can be described using one of the three following ways:

- English
- Pseudocode
- Programming language

Algorithm. (*Insertion Sort*)

Input : a_1, \dots, a_n , n unsorted elements

Output: the $a_i, 1 \leq i \leq n$, in increasing order

```
1 for  $j \leftarrow 2$  to  $n$  do
2    $i \leftarrow 1$ ;
3   while  $a_j > a_i$  do  $i \leftarrow i + 1$ ;
4    $m \leftarrow a_j$ ;
5   for  $k \leftarrow 0$  to  $j - i - 1$  do  $a_{j-k} \leftarrow a_{j-k-1}$ ;
6    $a_i \leftarrow m$ 
7 end for
8 return  $(a_1, \dots, a_n)$ 
```

Setup: a robot arm solders chips on a board in n contact points

Goal: minimize the time to attach a chip to the board

Assumptions:

- The arm moves at constant speed
- Once a chip has been attached another one is soldered

Setup: a robot arm solders chips on a board in n contact points

Goal: minimize the time to attach a chip to the board

Assumptions:

- The arm moves at constant speed
- Once a chip has been attached another one is soldered

Defining the Input and Output:

- Input: a set S of n points in the plane
- Output: the shortest path visiting all the points in S

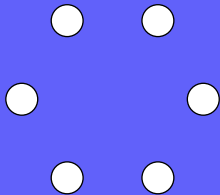
Algorithm. (*Nearest neighbor*)

Input : a set $S = \{s_0, \dots, s_{n-1}\}$ of n points in the plane

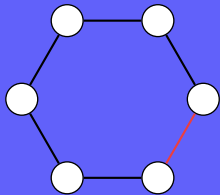
Output: the shortest cycle visiting all the points in S

```
1  $p_0 \leftarrow s_0$ ;  
2 for  $i \leftarrow 1$  to  $n - 1$  do  
3    $p_i \leftarrow$  closest unvisited neighbor to  $p_{i-1}$ ;  
4   Visit  $p_i$ ;  
5 end for  
6 return  $\langle p_0, \dots, p_{n-1} \rangle$ 
```

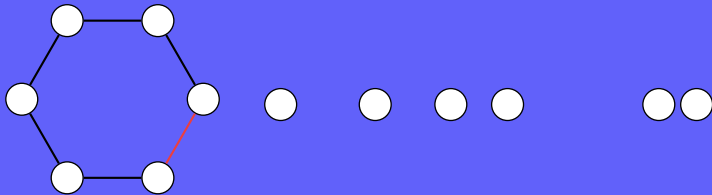
How does the nearest neighbor algorithm (1.26) perform in the following three cases?



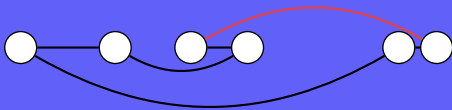
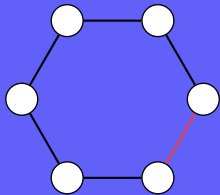
How does the nearest neighbor algorithm (1.26) perform in the following three cases?



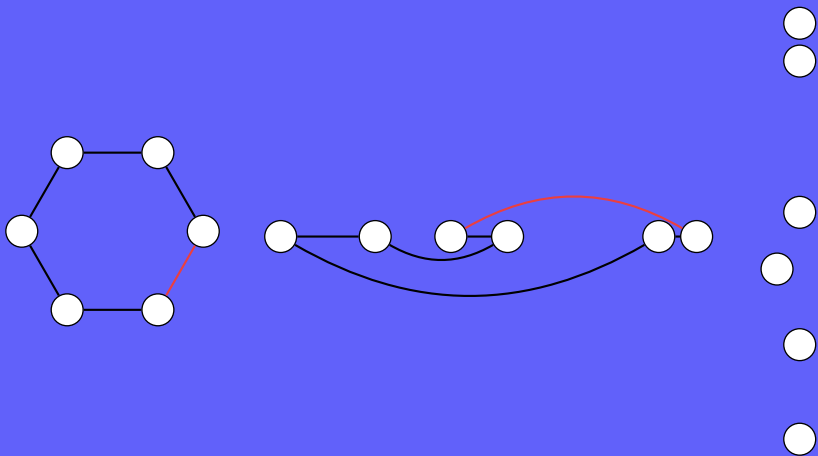
How does the nearest neighbor algorithm (1.26) perform in the following three cases?



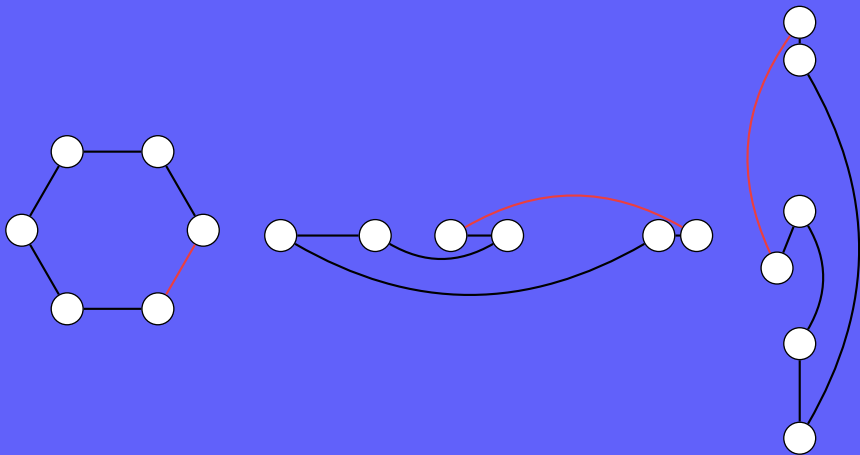
How does the nearest neighbor algorithm (1.26) perform in the following three cases?



How does the nearest neighbor algorithm (1.26) perform in the following three cases?



How does the nearest neighbor algorithm (1.26) perform in the following three cases?



Algorithm. (*Closest pair*)

Input : a set S of n points in the plane

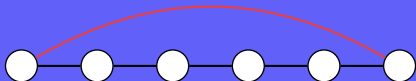
Output: the shortest cycle visiting all the points in S

```
1 for  $i \leftarrow 1$  to  $n - 1$  do
2    $d \leftarrow \infty$ ;
3   foreach pair of end points  $\langle s, t \rangle$  from distinct vertex chains do
4     if  $\text{dist}(s, t) \leq d$  then
5        $s_m \leftarrow s$ ;  $t_m \leftarrow t$ ;  $d \leftarrow \text{dist}(s, t)$ ;
6     end if
7   end foreach
8   Connect  $s_m$  and  $t_m$  by an edge;
9 end for
10 return all the points starting from one of the two end points
```

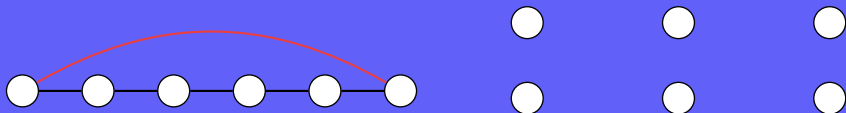
Applying the closest pair algorithm (1.28) to the following vertices arrangement yields the two graphs:



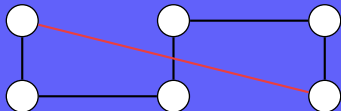
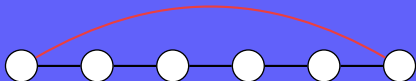
Applying the closest pair algorithm (1.28) to the following vertices arrangement yields the two graphs:



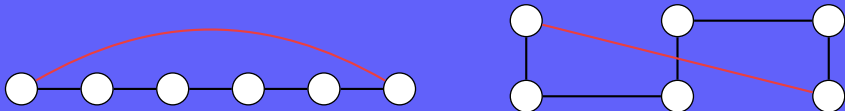
Applying the closest pair algorithm (1.28) to the following vertices arrangement yields the two graphs:



Applying the closest pair algorithm (1.28) to the following vertices arrangement yields the two graphs:



Applying the closest pair algorithm (1.28) to the following vertices arrangement yields the two graphs:



Possible strategy to ensure the most optimal path:

- Enumerate all the possible paths
- Select the one that minimizes the total length

Drawback: for only 20 vertices $20! = 2432902008176640000$ paths have to be explored...

A difference:

- Algorithm: always output a correct result
- Heuristic: idea serving as a guide to solve a problem with no guarantee of always providing a good solution

Correctness and efficiency:

- An algorithm working on a set of input does not imply it will work on all instances
- Efficient algorithm totally solving a problem might not exist

Common traps when defining the Input and Output:

- Are they precise enough?
- Can all the Input be easily and efficiently generated?
- Could there be any confusion on the expected Output?

Example.

For an Output, what does it mean to “find the best route”?

Common traps when defining the Input and Output:

- Are they precise enough?
- Can all the Input be easily and efficiently generated?
- Could there be any confusion on the expected Output?

Example.

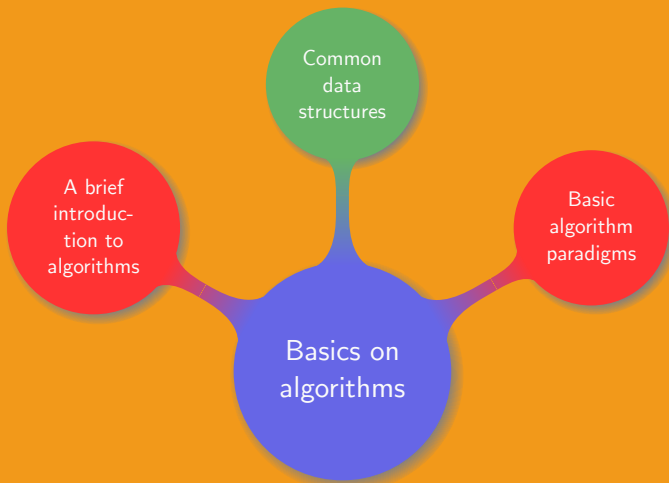
For an Output, what does it mean to “find the best route”?

The shortest in distance, the fastest in time, or the one minimizing the number of turns?

Conclusion: where to start (Input) and where to go (Output) must be expressed in simple, precise, and clear terms.

Finding good counter-examples:

- Seek simplicity: make it clear why the algorithm fails
- Think small: algorithms failing for large Input often fail for smaller one
- Test the extremes: study special cases, e.g. inputs equal, tiny, huge...
- Think exhaustively: test whether all the possible cases are covered by the algorithm
- Track weaknesses: check if the underlying idea behind the algorithm has any “unexpected” impact on the output



Data structures can be split into two main categories:

- **Continuous:** a single piece of memory
e.g. array, matrices, hash tables...
- **Linked data structures:** distinct chunks of memory connected together
e.g. linked list, trees, graph adjacency lists...

Each type has some relative advantages. Choosing the proper data type is of a major importance when designing algorithms.

Each element can be efficiently located using its index:

- **Constant access time:** each index maps to a memory address
- **Space efficiency:** no space wasted with links or information on the data
- **Memory locality:** data is contiguous so *cache* can be used to speed up successive data accesses

Downside: the size cannot be easily adjusted during the program's execution

A linked structure is composed of nodes. Each one contains:

- One or more fields on data
- A pointer to at least another node

The most common operations are:

- Search: find an item in the list
- Insert: add an item to the list
- Delete: remove an item from the list

Search can be implemented either iteratively or recursively

Linked structure

- Overflow only occurs when memory is full
- Insertion/deletion are simple and fast
- Moving pointers is faster than moving the actual data

Array

- No extra space wasted for the pointer field
- Efficient random access is possible
- Better memory locality and cache performance

Common data structures allowing the storage and retrieval of data independently of the content:

- **Stack:**
 - LIFO order
 - Simple to implement and very efficient
- **Queue:**
 - FIFO order
 - Minimize the maximum waiting time
 - Trickier to implement than stacks

Both can be implemented using either linked lists or arrays, depending if the size of the container is known in advance

Data type allowing access by content. Primary operations:

- **Search:** search a value in a given dictionary
- **Insert:** add an element to the dictionary
- **Delete:** remove an element from the dictionary

Most common operations:

- **Max/Min:** retrieve the largest/smallest element from the dictionary
- **Predecessor/Successor:** retrieve the element just before/after a given element; before/after are defined with respect to a sorted order

Let n be the number of elements in the array

Operation	Unsorted array	Sorted array
<code>search(D,k)</code>	$\mathcal{O}(n)$	$\mathcal{O}(\log n)$
<code>insert(D,k)</code>	$\mathcal{O}(1)$	$\mathcal{O}(n)$
<code>delete(D,k)</code>	$\mathcal{O}(1)^*$	$\mathcal{O}(n)$
<code>predecessor(D,k)</code>	$\mathcal{O}(n)$	$\mathcal{O}(1)$
<code>successor(D,k)</code>	$\mathcal{O}(n)$	$\mathcal{O}(1)$
<code>minimum(D)</code>	$\mathcal{O}(n)$	$\mathcal{O}(1)$
<code>maximum(D)</code>	$\mathcal{O}(n)$	$\mathcal{O}(1)$

* Assuming a pointer to the key k is given how to get $\mathcal{O}(1)$?

Dictionary using linked structures

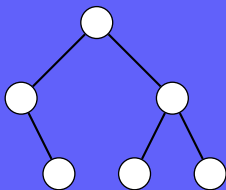
Let n be the number of elements in the list

Operation	Unsorted		Sorted	
	Single	Double	Single	Double
search(D,k)	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$
insert(D,k)	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$
delete(D,k)	$\mathcal{O}(n)^*$	$\mathcal{O}(1)$	$\mathcal{O}(n)^*$	$\mathcal{O}(1)$
predecessor(D,k)	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(n)^*$	$\mathcal{O}(1)$
successor(D,k)	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
minimum(D)	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
maximum(D)	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(1)^\dagger$	$\mathcal{O}(1)$

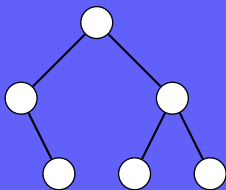
* Why are singly linked lists slower?

† How to achieve $\mathcal{O}(1)$ for singly sorted lists?

Binary search trees

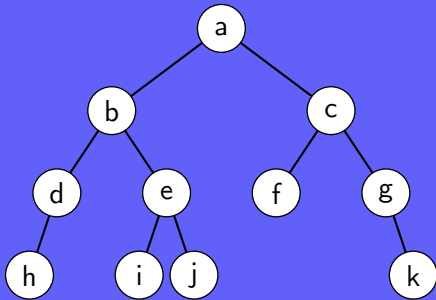


- Based on doubly linked lists
- First object is the root of the tree
- Second object is a left child if it precedes the root and a right child if it succeeds it
- Third and further object are sorted along the tree following a similar pattern



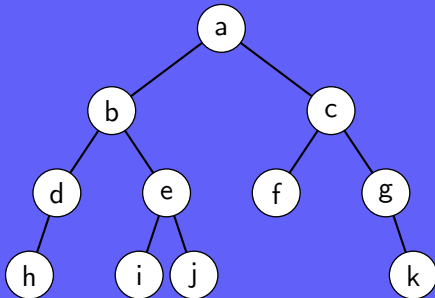
- Based on doubly linked lists
 - First object is the root of the tree
 - Second object is a left child if it precedes the root and a right child if it succeeds it
 - Third and further object are sorted along the tree following a similar pattern
-
- The three primary dictionary operations take $\mathcal{O}(h)$, with h the height of the tree
 - Binary search trees balance the search time and flexible update
 - How to handle deletion?

Binary search tree traversal



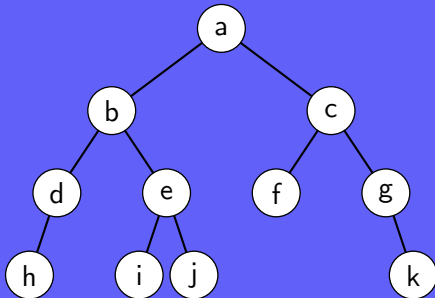
- Preorder traversal:

Binary search tree traversal



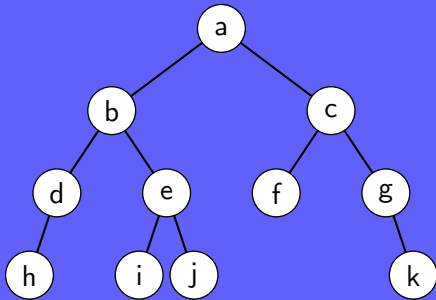
- Preorder traversal:
a, b, d, h, e, i, j, c, f, g, k
- Inorder traversal:

Binary search tree traversal



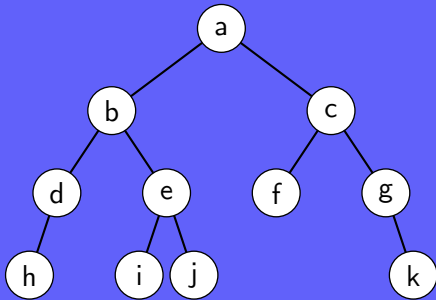
- Preorder traversal:
a, b, d, h, e, i, j, c, f, g, k
- Inorder traversal:
h, d, b, i, e, j, a, f, c, g, k
- Postorder traversal:

Binary search tree traversal



- Preorder traversal:
a, b, d, h, e, i, j, c, f, g, k
- Inorder traversal:
h, d, b, i, e, j, a, f, c, g, k
- Postorder traversal:
h, d, i, j, e, b, f, k, g, c, a

Binary search tree traversal



- Preorder traversal:
a, b, d, h, e, i, j, c, f, g, k
- Inorder traversal:
h, d, b, i, e, j, a, f, c, g, k
- Postorder traversal:
h, d, i, j, e, b, f, k, g, c, a

How to implement inorder tree traversal?

Primary operations for priority queues:

- **Insert:** add an element to the queue
- **Find min/max:** return the last/first element in the queue
- **Delete min/max:** remove the last/first element in the queue

Operation	Array		Balanced tree
	Unsorted	Sorted	
<code>insert(Q,x)</code>	$\mathcal{O}(1)$	$\mathcal{O}(n)$	$\mathcal{O}(\log n)$
<code>find_min(Q)</code>	$\mathcal{O}(1)^*$	$\mathcal{O}(1)$	$\mathcal{O}(1)^*$
<code>delete_min(Q)</code>	$\mathcal{O}(n)$	$\mathcal{O}(1)^\dagger$	$\mathcal{O}(\log n)$

* How to reach $\mathcal{O}(1)$ for an unsorted array and a balanced tree?

† How to reach $\mathcal{O}(1)$ when deleting the min in a sorted array?

Practical way to maintain a dictionary where:

- The data is stored in an array
- Each key is hashed and stored at index “the hash of the key”
- Keys with a similar hash are store in a linked list

Good hash function: all indices occur with equiprobability

Practical way to maintain a dictionary where:

- The data is stored in an array
- Each key is hashed and stored at index “the hash of the key”
- Keys with a similar hash are store in a linked list

Good hash function: all indices occur with equiprobability

Example.

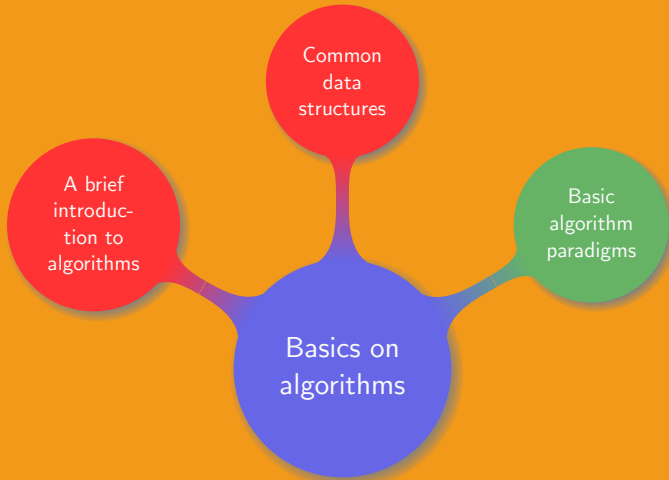
A common choice is $H(k) = k \bmod m$, with m a prime not too close from a power of 2.

For $n = 2000$, 701 would be a good choice if one desires to have about three keys stored at each index.

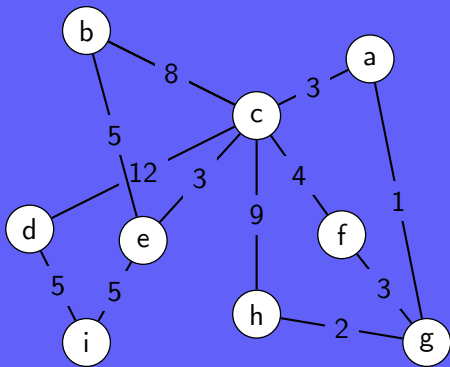
- **Strings:** array of characters; use suffix trees/arrays for pattern matching
- **Geometric element:** define regions as polygons using segments and points in an array or a tree
- **Graphs:** consider the adjacency matrix or an adjacency list; graph algorithms vary depending on the structure
- **Sets:** bit vector where the element in the set is the index and the value store is 1 or 0 depending whether the element is in the set; dictionaries can be used for fast membership queries

A few points to remember when selecting a data structure:

- Data can be represented in many ways
- No data structure is fast in all aspects
- Choosing the wrong data structure can be disastrous in terms of performance
- Several choices are often possible
- Identifying the best data structure is often not critical
- Always aim for clear, simple, and efficient data structures



A first graph problem



Simple problem:

- We have n computers connected by wires
- Using different wires implies different costs
- We want:
 - All the computers to be connected to the network
 - Minimize the cost

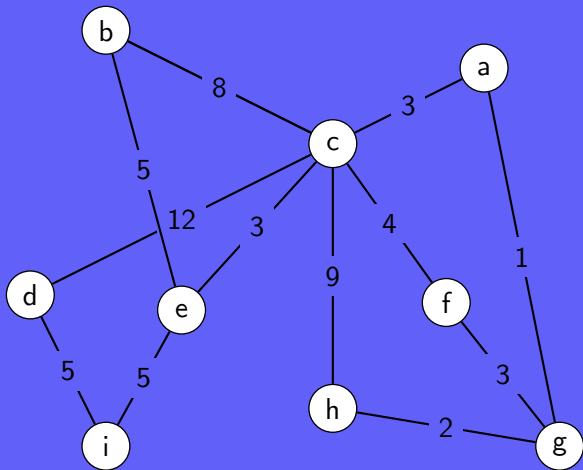
Problem (Minimum Spanning Tree (MST))

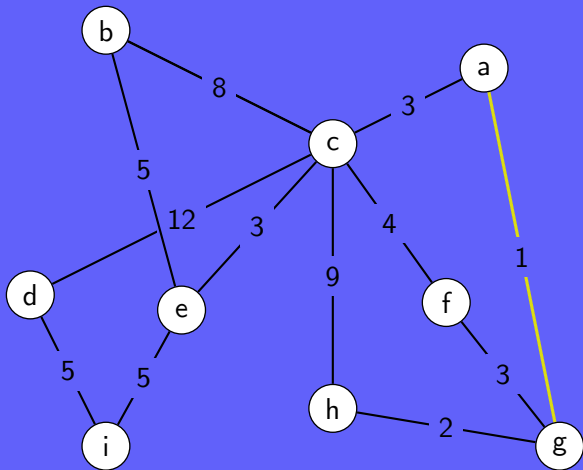
Given a weighted graph G , find a subgraph T such that:

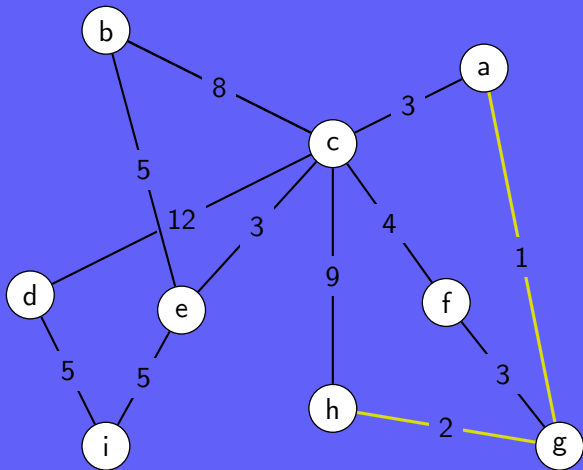
- ① All the vertices on G are connected on T ,
- ② The total weight, defined as the sum of the weight of all the edges in T , is minimized.

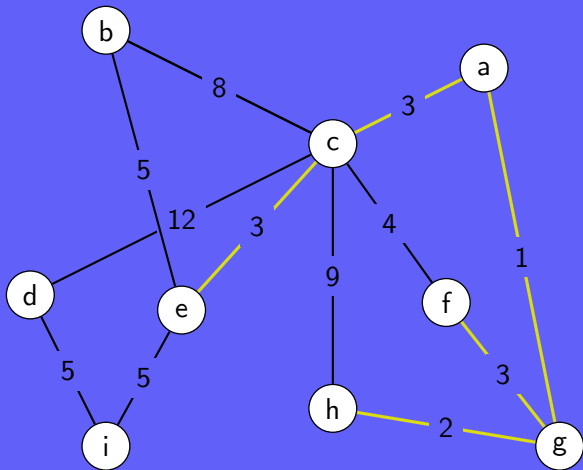
The graph T is a *minimum spanning tree* for G .

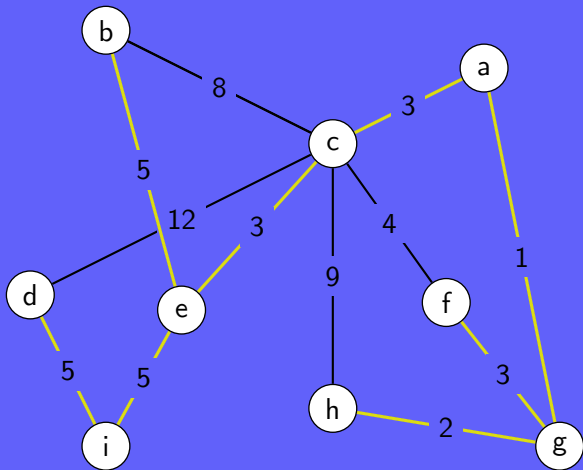
Note that T is a tree: if it contained a cycle, at least one edge could be removed, allowing a lower weight while preserving the connected property of T .











Algorithm. (*Kruskal*)

Input : A graph $G = \langle V, E \rangle$

Output: A minimum spanning tree T for G

```
1 Sort the edges  $G.E$  by weight;  
2  $T \leftarrow \emptyset$ ;  
3 for edges  $(u, v)$  in  $G.E$ , in non-decreasing order do  
4   | if adding  $(u, v)$  does not create a cycle then  
5   |   | add edge  $(u, v)$  to  $T$   
6   | end if  
7 end for  
8 return  $T$ 
```

Algorithm. (*Kruskal*)

Input : A graph $G = \langle V, E \rangle$

Output: A minimum spanning tree T for G

```
1 Sort the edges  $G.E$  by weight;  
2  $T \leftarrow \emptyset$ ;  
3 for edges  $(u, v)$  in  $G.E$ , in non-decreasing order do  
4   | if adding  $(u, v)$  does not create a cycle then  
5   |   | add edge  $(u, v)$  to  $T$   
6   | end if  
7 end for  
8 return  $T$ 
```

What needs to be specified?

Correctness of Kruskal's algorithm

Theorem

Assuming the previous notations, Kruskal's algorithm produces a minimum spanning tree for G .

Proof.

Let $G = \langle V, E \rangle$ be a graph and let v and w be two vertices connected by an edge. If S is the set of all the vertices with a path to v before e is added, then $w \notin S$, otherwise this would define a cycle. Moreover if there was an edge with smaller weight than e , connecting S and $V - S$, then it would have already been added. Therefore e is the cheapest edge connecting $V - S$ to S , and as such belongs to a minimum spanning tree of G .

Clearly by design the algorithm will not generate any cycle. Moreover as G is connected and all the edges are explored, $V - S$ and S will be linked at some stage. Hence T is connected. \square

Issue: how to represent the data such that whether or not adding an edge creates a cycle can be efficiently tested?

For each edge joining two vertices v and w :

- Identify all the connected components of v and w
- If the edge is to be included, merge the two components

Issue: how to represent the data such that whether or not adding an edge creates a cycle can be efficiently tested?

For each edge joining two vertices v and w :

- Identify all the connected components of v and w
- If the edge is to be included, merge the two components

Extra notes:

- No edge needs to be removed
- No component needs to be split
- Everything must be done efficiently

Representing data using:

- **An array:** testing can be done in constant time; merging requires linear time
- **A graph:** merging is only adding an edge; testing requires a full graph traversal

Representing data using:

- **An array:** testing can be done in constant time; merging requires linear time
- **A graph:** merging is only adding an edge; testing requires a full graph traversal

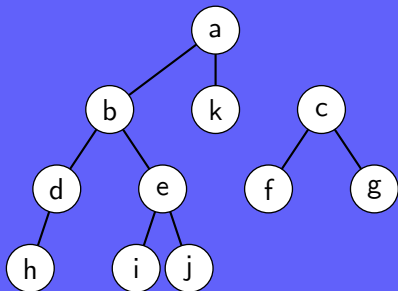
Implement a new data structure containing:

- A pointer to the parent
- The *rank*, or depth, of the sub-tree

The union-find data structure

The two operations are:

- $\text{Find}(v)$: find the root of the tree for vertex v
- $\text{Union}(v,w)$: link the root of the tree containing v to the root of the tree containing w (or the other way around)
- We have two sub-trees

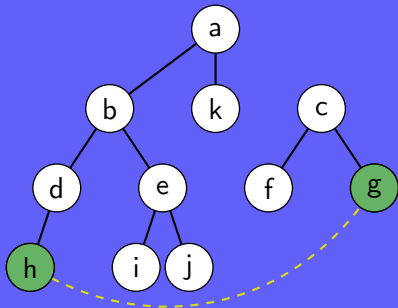


The union-find data structure

The two operations are:

- $\text{Find}(v)$: find the root of the tree for vertex v
- $\text{Union}(v, w)$: link the root of the tree containing v to the root of the tree containing w (or the other way around)

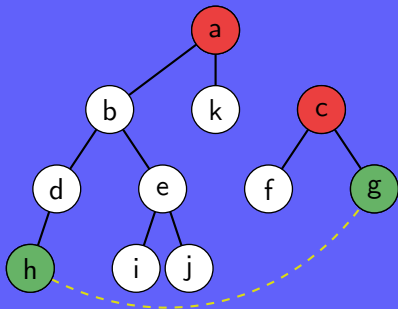
- We have two sub-trees
- On the graph an edge joins the vertices h and g



The union-find data structure

The two operations are:

- $\text{Find}(v)$: find the root of the tree for vertex v
- $\text{Union}(v, w)$: link the root of the tree containing v to the root of the tree containing w (or the other way around)

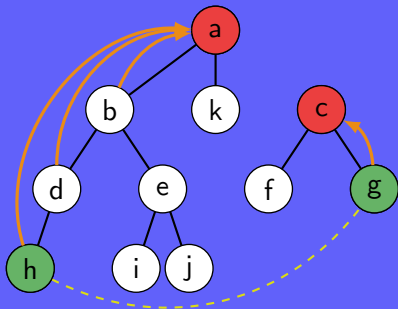


- We have two sub-trees
- On the graph an edge joins the vertices h and g
- Find on h and g returns a and c , respectively

The union-find data structure

The two operations are:

- $\text{Find}(v)$: find the root of the tree for vertex v
- $\text{Union}(v, w)$: link the root of the tree containing v to the root of the tree containing w (or the other way around)

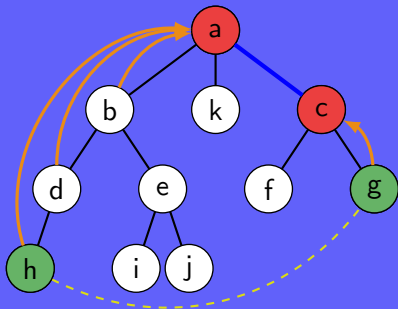


- We have two sub-trees
- On the graph an edge joins the vertices h and g
- Find on h and g returns a and c , respectively
- Update parents for h, d, b and g

The union-find data structure

The two operations are:

- $\text{Find}(v)$: find the root of the tree for vertex v
- $\text{Union}(v, w)$: link the root of the tree containing v to the root of the tree containing w (or the other way around)



- We have two sub-trees
- On the graph an edge joins the vertices h and g
- Find on h and g returns a and c , respectively
- Update parents for h, d, b and g
- Union connects c to a

Algorithm.

```
1 Function GenSet(x):  
2   |   x.parent  $\leftarrow$  x; x.rank  $\leftarrow$  0;  
3 end  
4 Function Find(x):  
5   |   if x.parent  $\neq$  x then x.parent  $\leftarrow$  Find(x.parent) ;  
6   |   return x.parent  
7 end  
8 Function Union(x,y):  
9   |   X  $\leftarrow$  Find(x); Y  $\leftarrow$  Find(y);  
10  |   if X.rank > Y.rank then Y.parent  $\leftarrow$  X;  
11  |   else X.parent  $\leftarrow$  Y;  
12  |   if X.rank = Y.rank then Y.rank++;  
13 end
```

Algorithm. (*Kruskal with find-union*)

Input : A graph $G = \langle V, E \rangle$

Output: A minimum spanning tree T

```
1 Sort the edges  $G.E$  by weight;
2  $T \leftarrow \emptyset$ ;
3 for edges  $(u, v)$  in  $G.E$ , in non-decreasing order do
4   | if Find( $u$ )  $\neq$  Find( $v$ ) then
5   |   | add edge  $(u, v)$  to  $T$ ;
6   |   | Union( $u, v$ )
7   | end if
8 end for
9 return  $T$ 
```

Problem (Counting inversions)

Given a list of n elements a_0, \dots, a_{n-1} , determine how many pairs $(a_i, a_j)_{\substack{0 \leq i, j \leq n \\ i \neq j}}$ are not in ascending order. Such pairs are called *inversions*.

Remark.

This problem has numerous applications such as

- Voting theory
- Analysis of search engines ranking
- Collaborative filtering

Given 6 movies compare the ranking of two users:

Movie	A	B	C	D	E	F
First user	1	2	3	4	5	6
Second user	1	3	5	2	4	6

Given 6 movies compare the ranking of two users:

Movie	A	B	C	D	E	F
First user	1	2	3	4	5	6
Second user	1	3	5	2	4	6

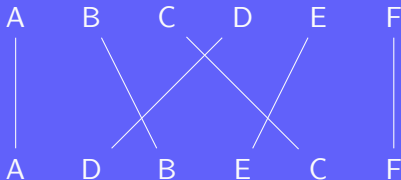
Inversions: $(3, 2)$, $(5, 2)$, $(5, 4)$

Given 6 movies compare the ranking of two users:

Movie	A	B	C	D	E	F
First user	1	2	3	4	5	6
Second user	1	3	5	2	4	6

Inversions: $(3, 2)$, $(5, 2)$, $(5, 4)$

A simple geometrical view:



Divide and conquer approach

Strategy for solving the counting inversions problem (1.59):

- 1 **Divide:** split the list L into two halves L_1 and L_2
- 2 **Conquer:** recursively count inversions in each list
- 3 **Combine:** count inversions for the pairs (l_i, l_j) with l_i and l_j belonging to L_1 and L_2 respectively

The sum of the three counts is the total number of inversion in L

Example.

1 5 4 8 10 2 6 9 3 7

1 5 4 8 10

2 6 9 3 7

Divide and conquer approach

Strategy for solving the counting inversions problem (1.59):

- 1 **Divide:** split the list L into two halves L_1 and L_2
- 2 **Conquer:** recursively count inversions in each list
- 3 **Combine:** count inversions for the pairs (l_i, l_j) with l_i and l_j belonging to L_1 and L_2 respectively

The sum of the three counts is the total number of inversion in L

Example.

1 5 4 8 10 2 6 9 3 7

1 5 4 8 10

(5,4)

2 6 9 3 7

(6,3),(9,3),(9,7)

1 4 5 8 10

2 3 6 7 9

Divide and conquer approach

Strategy for solving the counting inversions problem (1.59):

- ① **Divide:** split the list L into two halves L_1 and L_2
- ② **Conquer:** recursively count inversions in each list
- ③ **Combine:** count inversions for the pairs (l_i, l_j) with l_i and l_j belonging to L_1 and L_2 respectively

The sum of the three counts is the total number of inversion in L

Example.

1 5 4 8 10 2 6 9 3 7

1 5 4 8 10

(5,4)

2 6 9 3 7

(6,3),(9,3),(9,7)

1 4 5 8 10

2 3 6 7 9

(4,2),(4,3),(5,2),(5,3),(8,2),(8,3),(8,6),(8,7),(10,2),(10,3),(10,6),(10,7),(10,9)

Algorithm. (*Merge and count*)

Input : Two sorted lists: $L_1 = (l_{1,1}, \dots, l_{1,n_1})$, $L_2 = (l_{2,1}, \dots, l_{2,n_2})$

Output: Number of inversions *count*, and L_1 and L_2 merged into L

```
1 Function MergeCount( $L_1, L_2$ ):  
2    $count \leftarrow 0$ ;  $L \leftarrow \emptyset$ ;  $i \leftarrow 1$ ;  $j \leftarrow 1$ ;  
3   while  $i \leq n_1$  and  $j \leq n_2$  do  
4     if  $l_{1,i} \leq l_{2,j}$  then  
5        $\text{append } l_{1,i} \text{ to } L$ ;  $i++$ ;  
6     else  
7        $\text{append } l_{2,j} \text{ to } L$ ;  $count \leftarrow count + n_1 - i + 1$ ;  $j++$ ;  
8     end if  
9   end while  
10  if  $i > n_1$  then  $\text{append } l_{2,j}, \dots, l_{2,n_2} \text{ to } L$ ;  
11  else  $\text{append } l_{1,i}, \dots, l_{1,n_1} \text{ to } L$ ;  
12  return  $count$  and  $L$   
13 end
```

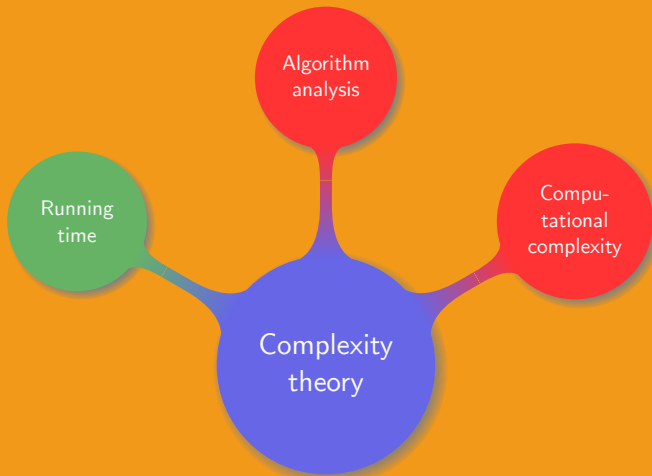
Algorithm. (*Sort and count*)

Input : A list $L = (l_1, \dots, l_n)$

Output: The number of inversions *count* and L sorted

```
1 Function SortCount( $L$ ):  
2   if  $n=1$  then return 0 and  $L$ ;  
3   else  
4     Split  $L$  into  $L_1 = (l_1, \dots, l_{\lceil n/2 \rceil})$  and  $L_2 = (l_{\lceil n/2 \rceil + 1}, \dots, l_n)$ ;  
5      $count_1, L_1 \leftarrow$  SortCount( $L_1$ );  
6      $count_2, L_2 \leftarrow$  SortCount( $L_2$ );  
7      $count, L \leftarrow$  MergeCount( $L_1, L_2$ );  
8   end if  
9    $count \leftarrow count_1 + count_2 + count$ ;  
10  return  $count$  and  $L$   
11 end
```

2. Complexity theory



Random Access Machine (RAM) model:

- Each simple operation ($+$, $-$, \times , $/$, if, \dots) takes one step
- Loops are composed of several simple steps, which are repeated a certain number of times
- Each memory access accounts for one step

The runtime of an algorithm is given by the total number of steps necessary to complete it.

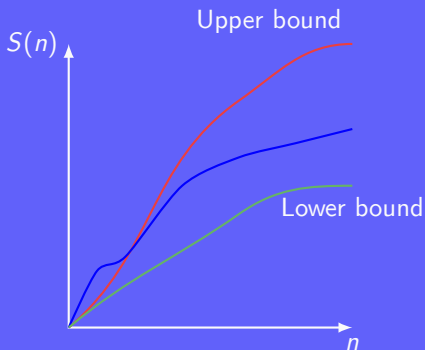
The RAM model allows the study of algorithms independently of their implementation or running environment.



Given a problem with input size n :

- **Worst case complexity:** maximum number of steps to complete an instance of the algorithm
- **Best case complexity:** minimum number of steps to complete an instance of the algorithm
- **Average case complexity:** average number of steps, over all possible instances

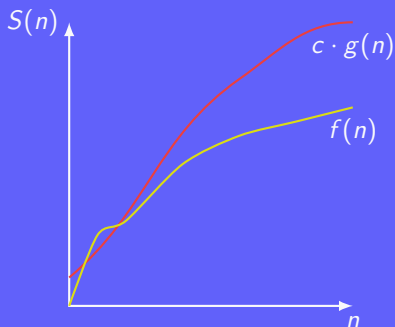
The complexity of an algorithm is defined by a numerical function.



Limitations of the RAM:

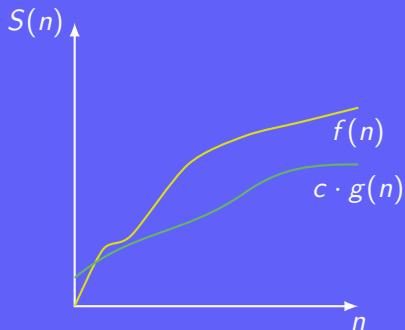
- Complexity can display many bumps
e.g. algorithms performing better on input of size a power of 2
- Counting the exact number of operations is too complex
e.g. implementation choices impact the number of RAM instructions

We need notations that simplify the analysis by ignoring the details irrelevant to the comparison of the algorithms.



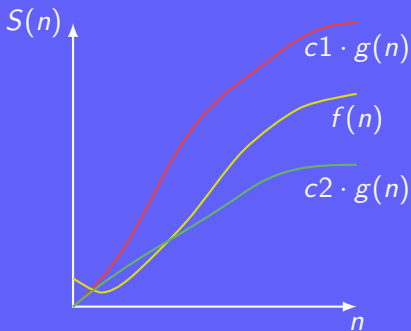
Approximating complexity:

- $c \cdot g(n)$ is an upper bound for $f(n)$
- There exists two constants c and n_0 such that for all n larger than n_0 ,
 $f(n) \leq c \cdot g(n)$
- $f(n) = \mathcal{O}(g(n))$



Approximating complexity:

- $c \cdot g(n)$ is a lower bound for $f(n)$
- There exists two constants c and n_0 such that for all n larger than n_0 , $f(n) \geq c \cdot g(n)$
- $f(n) = \Omega(g(n))$



Approximating complexity:

- $c_1 \cdot g(n)$ is an upper bound for $f(n)$ while $c_2 \cdot g(n)$ is a lower bound for $f(n)$
- There exists three constants c_1 , c_2 , and n_0 such that for all n larger than n_0 , $c_2 \cdot g(n) \leq f(n) \leq c_1 \cdot g(n)$
- $f(n) = \Theta(g(n))$

Exercise.

Prove the following equalities:

- $2^{n+1} = \Theta(2^n)$?

Exercise.

Prove the following equalities:

- $2^{n+1} = \Theta(2^n)$?

It suffices to prove that $2^{n+1} = \mathcal{O}(2^n)$ and $2^{n+1} = \Omega(2^n)$.

First note that $2^{n+1} = 2 \cdot 2^n$, such that $2^{n+1} \leq c \cdot 2^n$, for any $c \geq 2$. Then observe that for all $n \geq 2$, $2^{n+1} \geq c \cdot 2^n$, for any $0 < c \leq 2$.

- $(x + y)^2 = \mathcal{O}(x^2 + y^2)$?

Exercise.

Prove the following equalities:

- $2^{n+1} = \Theta(2^n)$?

It suffices to prove that $2^{n+1} = \mathcal{O}(2^n)$ and $2^{n+1} = \Omega(2^n)$.

First note that $2^{n+1} = 2 \cdot 2^n$, such that $2^{n+1} \leq c \cdot 2^n$, for any $c \geq 2$. Then observe that for all $n \geq 2$, $2^{n+1} \geq c \cdot 2^n$, for any $0 < c \leq 2$.

- $(x + y)^2 = \mathcal{O}(x^2 + y^2)$?

We need to find some c such that $(x + y)^2 \leq c(x^2 + y^2)$.

Without loss of generality we can assume $x \geq y$. Then

$2xy \leq 2x^2 \leq 2(x^2 + y^2)$ and as $(x + y)^2 = x^2 + 2xy + y^2$, we conclude that $(x + y)^2 \leq 3(x^2 + y^2)$.

Common complexity functions:

- Constant: $\mathcal{O}(1)$
- Logarithmic: $\mathcal{O}(\log n)$
- Linear: $\mathcal{O}(n)$
- Superlinear: $\mathcal{O}(n \log n)$
- Quadratic: $\mathcal{O}(n^2)$
- Cubic: $\mathcal{O}(n^3)$
- Exponential: $\mathcal{O}(c^n)$
- Factorial: $\mathcal{O}(n!)$

Common complexity functions:

- Constant: $\mathcal{O}(1)$
- Logarithmic: $\mathcal{O}(\log n)$
- Linear: $\mathcal{O}(n)$
- Superlinear: $\mathcal{O}(n \log n)$
- Quadratic: $\mathcal{O}(n^2)$
- Cubic: $\mathcal{O}(n^3)$
- Exponential: $\mathcal{O}(c^n)$
- Factorial: $\mathcal{O}(n!)$

More complexity functions:

- Inverse Ackerman's function: $\mathcal{O}(\alpha(n))$
- Log Log: $\mathcal{O}(\log \log n)$
- $\mathcal{O}(\log n / \log \log n)$
- $\mathcal{O}(\sqrt{n})$
- $\mathcal{O}(n^{c+\epsilon})$

Common complexity functions:

- Constant: $\mathcal{O}(1)$
- Logarithmic: $\mathcal{O}(\log n)$
- Linear: $\mathcal{O}(n)$
- Superlinear: $\mathcal{O}(n \log n)$
- Quadratic: $\mathcal{O}(n^2)$
- Cubic: $\mathcal{O}(n^3)$
- Exponential: $\mathcal{O}(c^n)$
- Factorial: $\mathcal{O}(n!)$

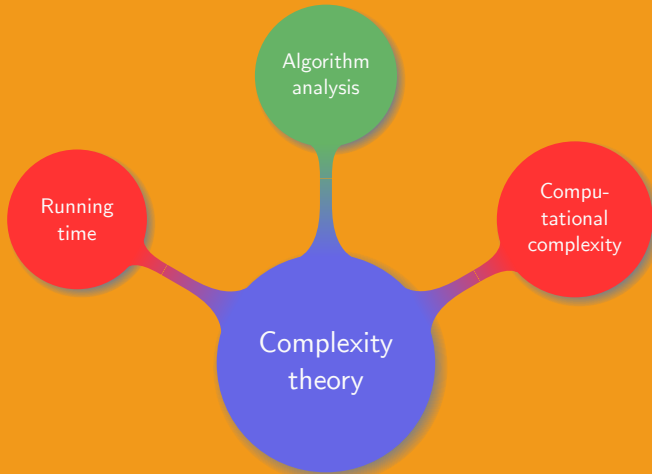
More complexity functions:

- Inverse Ackerman's function: $\mathcal{O}(\alpha(n))$
- Log Log: $\mathcal{O}(\log \log n)$
- $\mathcal{O}(\log n / \log \log n)$
- $\mathcal{O}(\sqrt{n})$
- $\mathcal{O}(n^{c+\epsilon})$

$$\begin{aligned} n! &\gg 2^n \gg n^{c+\epsilon} \gg n^3 \gg n^{2+\epsilon} \gg n^2 \gg n \log n \gg n^{1+\epsilon} \gg n \\ &\gg \sqrt{n} \gg \log n \gg \frac{\log n}{\log \log n} \gg \log \log n \gg \alpha(n) \gg 1 \end{aligned}$$

Assuming 10^9 basic operations a second, the time spent depending on the input size (in bits or number of input) and the time complexity is given as follows.

Input size n	$\log n$	n	$n \log n$	n^2	n^3	2^n	$n!$
10	$< 0.01\mu s$	$0.01\mu s$	$0.03\mu s$	$0.1\mu s$	$1\mu s$	$1\mu s$	3.6ms
20	$< 0.01\mu s$	$0.02\mu s$	$0.09\mu s$	$0.4\mu s$	$0.8\mu s$	1ms	77y
30	$< 0.01\mu s$	$0.03\mu s$	$0.1\mu s$	$0.9\mu s$	$2.7\mu s$	1s	$10^{16}y$
40	$< 0.01\mu s$	$0.04\mu s$	$0.2\mu s$	$1.6\mu s$	$6.4\mu s$	18min	
50	$< 0.01\mu s$	$0.05\mu s$	$0.3\mu s$	$2.5\mu s$	$12.5\mu s$	13d	
100	$< 0.01\mu s$	$0.1\mu s$	$0.6\mu s$	$10\mu s$	1ms	$10^{13}y$	
1000	$0.01\mu s$	$1\mu s$	$10\mu s$	1ms	1s		
10000	$0.01\mu s$	$10\mu s$	$130\mu s$	100ms	16min		
100000	$0.02\mu s$	$100\mu s$	$1.7ms$	10s	11.5d		
1000000	$0.02\mu s$	1ms	20ms	17min	32y		
10000000	$0.02\mu s$	10ms	200ms	1.2d	32000y		
100000000	$0.03\mu s$	100ms	2.6s	116d			
1000000000	$0.03\mu s$	1s	30s	32y			



Stable marriage problem

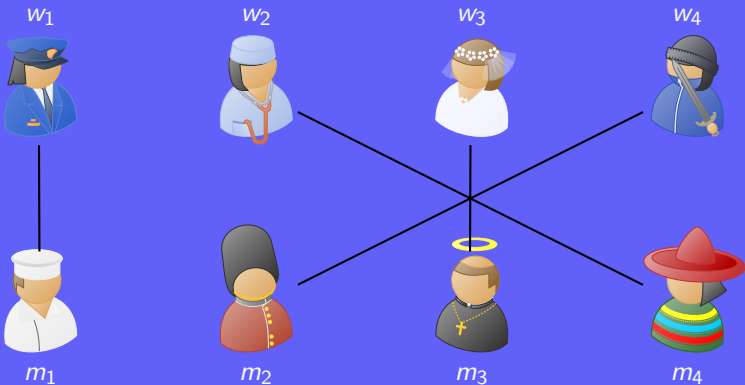
 w_1  w_2  w_3  w_4  m_1  m_2  m_3  m_4

Stable marriage problem

 $m_3 > m_2 > m_1 > m_4$  $m_4 > m_1 > m_2 > m_3$  $m_3 > m_4 > m_1 > m_2$  $m_1 > m_2 > m_3 > m_4$  $w_1 > w_2 > w_3 > w_4$ $w_1 > w_4 > w_3 > w_2$ $w_1 > w_3 > w_4 > w_2$ $w_1 > w_3 > w_2 > w_4$

Stable marriage problem

$m_3 > m_2 > m_1 > m_4$
 $m_4 > m_1 > m_2 > m_3$
 $m_3 > m_4 > m_1 > m_2$
 $m_1 > m_2 > m_3 > m_4$



$w_1 > w_2 > w_3 > w_4$

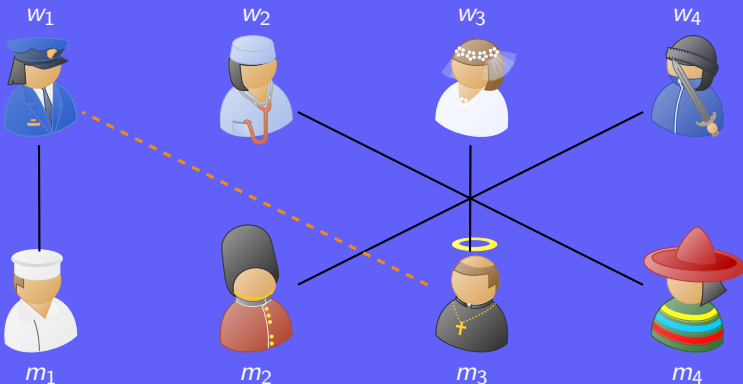
$w_1 > w_4 > w_3 > w_2$

$w_1 > w_3 > w_4 > w_2$

$w_1 > w_3 > w_2 > w_4$

Stable marriage problem

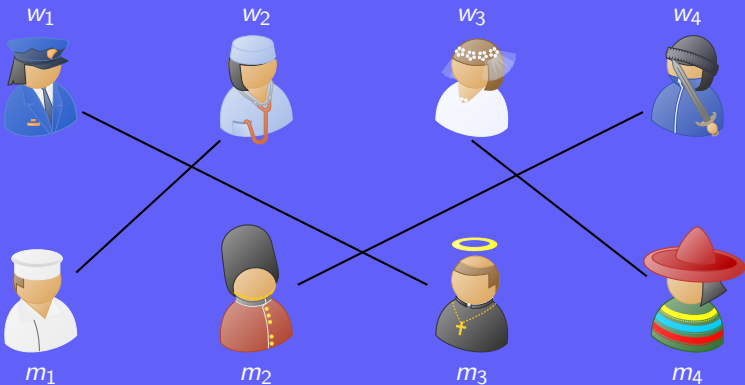
$m_3 > m_2 > m_1 > m_4$
 $m_4 > m_1 > m_2 > m_3$
 $m_3 > m_4 > m_1 > m_2$
 $m_1 > m_2 > m_3 > m_4$



$w_1 > w_2 > w_3 > w_4$
 $w_1 > w_4 > w_3 > w_2$
 $w_1 > w_3 > w_4 > w_2$
 $w_1 > w_3 > w_2 > w_4$

Stable marriage problem

$m_3 > m_2 > m_1 > m_4$
 $m_4 > m_1 > m_2 > m_3$
 $m_3 > m_4 > m_1 > m_2$
 $m_1 > m_2 > m_3 > m_4$



$w_1 > w_2 > w_3 > w_4$

$w_1 > w_4 > w_3 > w_2$

$w_1 > w_3 > w_4 > w_2$

$w_1 > w_3 > w_2 > w_4$

Problem (Stable matching problem)

Given a set of n men and a set n women, each person ranks all the members of the other set. A match between the men and women where no $(man, woman)$ -pair, would prefer to be together rather than with their current partner is said to be *stable*.

Problem (Stable matching problem)

Given a set of n men and a set n women, each person ranks all the members of the other set. A match between the men and women where no $(man, woman)$ -pair, would prefer to be together rather than with their current partner is said to be *stable*.

Visual representation of the problem:

- Each men (women) is represented by a vertex
- No edge connects elements from a same set
- Monogamy is assumed for both men and women
- The number of edges equals the number of vertices in a set
- No blocking pair must exist

Algorithm. (*Gale-Shapley*)

Input : n men and n women, all initially free

Output : n engaged pairs

```
1 while there is a free man do
2    $m \leftarrow$  select a man;  $w \leftarrow$  favorite woman  $m$  hasn't proposed yet;
3    $m$  proposes to  $w$ ;
4   if  $w$  is free then set  $(m, w)$  as an engaged pair;
5   else
6      $m' \leftarrow$  current date of  $w$ ;
7     if  $w$  prefers  $m'$  over  $m$  then  $m$  remains free;
8     else
9       set  $(m, w)$  as an engaged pair;
10      set  $m'$  as free;
11    end if
12  end if
13 end while
14 return the  $n$  engaged pairs
```

Theorem

Given two sets of n men and n women, all initially free, Gale-Shapley algorithm returns a stable matching.

Proof.

First we prove that algorithm 2.79 returns a perfect matching, i.e. there is no free man or woman.

Suppose there exists a free man who has already proposed to all the women. As a woman who is engaged once remains engaged, to the same man or a new one, it means that all the n women are engaged. Noting that a woman cannot be engaged to more than one man leads to a contradiction. ⚡

Lets now assume that the matching is not stable. Then there are two pairs (m, w) and (m', w') such that m prefers w' over w while w' prefers m over m' .

Proof (Continued).

By definition, the last proposal of m was to w . If he didn't proposed earlier to w' , then he does not prefer w' over w . ⚡

If he did propose to w' , then he was rejected in favor of another man m'' . Therefore the final partner of w' is either m' or a man she prefers over m' . Either way w' does not prefer m to m' . ⚡

It follows that Gale-Shapley algorithm returns a stable matching.



Proof (Continued).

By definition, the last proposal of m was to w . If he didn't proposed earlier to w' , then he does not prefer w' over w . ⚡

If he did propose to w' , then he was rejected in favor of another man m'' . Therefore the final partner of w' is either m' or a man she prefers over m' . Either way w' does not prefer m to m' . ⚡

It follows that Gale-Shapley algorithm returns a stable matching.



Corollary

If each man prefers a different woman then they all end up with their first choice, independently of the women preferences.

Theorem

Gale-Shapley algorithm has complexity $\mathcal{O}(n^2)$.

Proof.

At each iteration of the `while` loop a man proposes to a woman he has never proposed before. Therefore, for n men, there is a maximum of n^2 proposals.

Since each iteration corresponds to exactly one proposal the `while` loop is applied at most n^2 times.

Hence algorithm 2.79 has complexity $\mathcal{O}(n^2)$.



Complexity of the Union-Find structure

In order to study the complexity of the Union-Find data structure (1.56) we introduce the following simple results.

Lemma

- ① A node that is a root and gets attached to another root will never be a root again.
- ② When a node stops being a root its rank remains fixed.
- ③ As Find travels to the root of the tree the rank of the nodes strictly increases.
- ④ A node with rank k has at least 2^k nodes in its subtree.
- ⑤ Over n elements there are at most $n/2^k$ elements of rank k .

Definition

The *iterated logarithm* function, denoted \log^* , is defined by

$$\log^* : \mathbb{R} \longrightarrow \mathbb{N}$$
$$x \longmapsto \begin{cases} 0 & \text{if } x \leq 1 \\ 1 + \log^* \log_2 x & \text{if } x > 1 \end{cases}$$

In a less abstract way the iterated logarithm of n is the number of times the logarithm function has to be applied in order to get a number smaller than 2.

Definition

The *iterated logarithm* function, denoted \log^* , is defined by

$$\log^* : \mathbb{R} \longrightarrow \mathbb{N}$$
$$x \longmapsto \begin{cases} 0 & \text{if } x \leq 1 \\ 1 + \log^* \log_2 x & \text{if } x > 1 \end{cases}$$

In a less abstract way the iterated logarithm of n is the number of times the logarithm function has to be applied in order to get a number smaller than 2.

Example.

$$\log^* 4 = 1 + \log^* \log_2 4 = 2$$

$$\log^* 16 = 1 + \log^* 4 = 3$$

$$\log^* 65536 = \log^* 2^{2^{2^2}} = 4$$

Theorem

The cost of one Find in the Union-Find data structure is $\mathcal{O}(\log n)$, while m Find operations cost $\mathcal{O}((m + n) \log^* n)$.

Proof.

The rank of the root being bounded by the depth of the tree in the Union-Find data structure it is at most $\log n$.

The problem is now to evaluate the effect of m such operations. Since the path is compressed at each Find the complexity of any subsequent Find reusing a visited path is decreased.

In fact the running time of a Find operation is proportional to (i) the length of the path from a node to the root of the tree it belongs to and (ii) m .

Proof (continued).

We now divide the nodes into “blocks”, such that nodes of similar ranks are on a same block. To formally define a block we introduce the following recursive sequence

$$\begin{cases} T_0 = 1 \\ T_i = 2^{T_{i-1}}, \quad i > 0. \end{cases}$$

The blocks are then defined as the sets

$$\begin{cases} B_0 = \{0\}, \quad B_1 = \{1\}, \\ B_i = [T_{i-1}, T_i - 1], \quad i > 1. \end{cases}$$

Observing the blocks we notice that the maximum number of blocks is $\log^* n$.

Proof (continued).

Moreover as there is a maximum of $n/2^k$ elements of rank k (lemma 2.83), the number of nodes in the i -th block is at most

$$\begin{aligned}\sum_{j=T_{i-1}}^{T_i-1} \frac{n}{2^j} &= n \sum_{j=T_{i-1}}^{2^{T_{i-1}}-1} \frac{1}{2^j} \\ &\leq \frac{2n}{2^{T_{i-1}}} \\ &= \frac{2n}{T_i}.\end{aligned}\tag{2.1}$$

The idea is to evaluate the cost of the `Find` while traveling from a node to its root. In such a case we traverse nodes which are (i) in the same block (ii) in a different block or (iii) directly connected to the root.

Proof (continued.)

The simplest case is (iii) as only one step is required, implying an $\mathcal{O}(m)$ complexity for the m Find.

Note that (ii) was evaluated earlier when we observed that the maximum number of blocks is $\log^* n$. Therefore the complexity of (ii) is $\mathcal{O}(m \log^* n)$.

Finally for case (i) the path starting at a node x goes through at most $T_i - 1 - T_{i-1}$ ranks, and from (2.1) there are less than $2n/T_i$ such elements x in a block. Thus the total number of times the rank changes in a block is

$$\begin{aligned} (T_i - 1 - T_{i-1}) \frac{2n}{T_i} &\leq T_i \frac{2n}{T_i} \\ &= 2n. \end{aligned}$$

Proof (continued).

Remembering that the maximum number of blocks is $\log^* n$, the Find operations generate an overall of at most $\mathcal{O}(n \log^* n)$ internal links.

Hence the time spent on m operations is given by the total time spent on cases (i) to (iii), that is $\mathcal{O}((m + n) \log^* n)$. \square

For long this has been the tightest proven upper bound. However a better result was proven in the 1970's. It bounds the complexity by using the inverse Ackerman's function, which is growing even slower than the \log^* function.

Definition (Ackerman's function)

For any two positive integers m and n , Ackerman's function is recursively defined by

$$A(m, n) = \begin{cases} n + 1 & \text{if } m = 0 \\ A(m - 1, 1) & \text{if } m > 0 \text{ and } n = 0 \\ A(m - 1, A(m, n - 1)) & \text{if } m > 0 \text{ and } n > 0 \end{cases}$$

Definition (Ackerman's function)

For any two positive integers m and n , Ackerman's function is recursively defined by

$$A(m, n) = \begin{cases} n + 1 & \text{if } m = 0 \\ A(m - 1, 1) & \text{if } m > 0 \text{ and } n = 0 \\ A(m - 1, A(m, n - 1)) & \text{if } m > 0 \text{ and } n > 0 \end{cases}$$

Example.

$$A(0, 0) = 1$$

$$A(1, 1) = A(0, A(1, 0)) = A(0, 2) = 3$$

$$A(2, 2) = A(1, A(2, 1)) = A(1, 5) = 7$$

$$A(3, 3) = A(2, A(3, 2)) = A(2, 29) = 61$$

$$A(4, 4) = A(3, A(4, 3)) = 2^{2^{2^{65536}}} - 3$$

Inverse Ackerman's function

For $m = n$ one defines the inverse of Ackerman's function $\alpha(x)$. As Ackerman's function is extremely fast-growing its inverse is an extremely slow growing function which never increases beyond 4 or 5 for any practical value.

Theorem

The amortized time for a sequence of m GenSet, Union, and Find operations, n of which are GenSet, can be performed in time $\mathcal{O}(m\alpha(n))$, where $\alpha(n)$ is the inverse Ackerman's function.

Inverse Ackerman's function

For $m = n$ one defines the inverse of Ackerman's function $\alpha(x)$. As Ackerman's function is extremely fast-growing its inverse is an extremely-slow growing function which never increases beyond 4 or 5 for any practical value.

Theorem

The amortized time for a sequence of m `GenSet`, `Union`, and `Find` operations, n of which are `GenSet`, can be performed in time $\mathcal{O}(m\alpha(n))$, where $\alpha(n)$ is the inverse Ackerman's function.

Remark.

It was recently proven that this upper bound is asymptotically optimal, that is the complexity of the Union-Find structure is $\Omega(\alpha(n))$.

The Sort and Count algorithm (1.63) solves the Counting inversions problem (1.59).

For the merge part:

- Each iteration of the `While` takes constant time
- Elements that are added are never added again
- There is a maximum of n element

The time spent for the merge part is $\mathcal{O}(n)$.

The Sort and Count algorithm (1.63) solves the Counting inversions problem (1.59).

For the merge part:

- Each iteration of the `While` takes constant time
- Elements that are added are never added again
- There is a maximum of n element

The time spent for the merge part is $\mathcal{O}(n)$.

For the sort part the divide and conquer strategy is applied:

- Divide the input into two pieces of equal size
- Recursively solve the two separate problems
- Combine the two results

The time spent for the division and recombination processes is linear in the size of the input.

Let $S(n)$ be the worst-case running time for sorting an instance of size n .

Simple complexity analysis:

- Time spent to divide into two pieces: $\mathcal{O}(n)$
- Time spent to solve each piece: $S(n/2)$
- Combine the results: $\mathcal{O}(n)$

Let $S(n)$ be the worst-case running time for sorting an instance of size n .

Simple complexity analysis:

- Time spent to divide into two pieces: $\mathcal{O}(n)$
- Time spent to solve each piece: $S(n/2)$
- Combine the results: $\mathcal{O}(n)$

Conclusion: the running time satisfies the recurrence relation

$$\begin{cases} S(2) \leq c & \text{for some constant } c \\ S(n) \leq 2S(n/2) + cn & \text{if } n > 2 \end{cases}$$

Let A be an algorithm whose running time is described by a recurrence of the form $T(n) = aT(n/b) + f(n)$, with $a \geq 1$, $b > 1$ two constants, and $f(n)$ an asymptotically positive function. This relation corresponds to A dividing the initial problem of size n into a sub-problems of size n/b each. While it takes $T(n/b)$ to recursively solve each of the a sub-problems $f(n)$ corresponds to the time spent splitting them and combining their results.

Remark.

For the sake of simplicity we assume n/b to be an integer, and more generally n to be a power of b . Although it is often not true, this simplifies the discussion while having no impact on the asymptotic behavior of the recurrence.

Theorem (Master theorem)

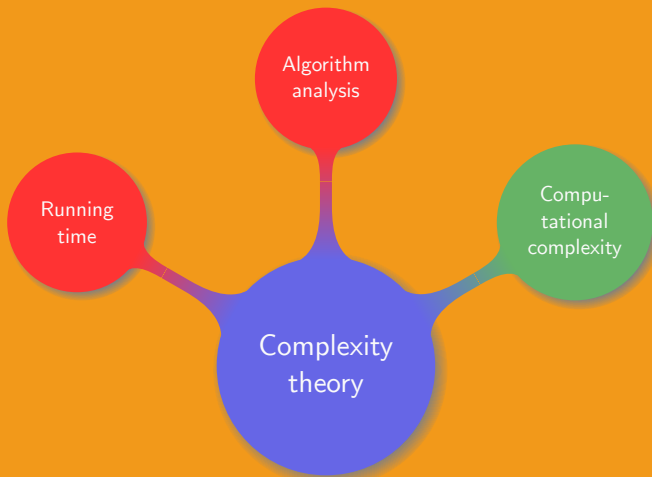
Let $a \geq 1$, $b > 1$, be two constants, $f(n)$ be a function, and $T(n) = aT(n/b) + f(n)$ be a recurrence relation over the positive integers. Then the asymptotic bound on $T(n)$ is given by

$$T(n) = \begin{cases} \Theta(n^{\log_b a} \log n) & \text{if } f(n) = \Theta(n^{\log_b a}); \\ \Theta(n^{\log_b a}) & \text{if } f(n) = \mathcal{O}(n^{\log_b a - \epsilon}), \epsilon > 0; \\ \Theta(f(n)) & \text{if } f(n) = \Omega(n^{\log_b a + \epsilon}), \epsilon > 0, n \text{ large} \\ & \text{enough, and } af(n/b) \leq cf(n), c < 1; \end{cases}$$

Summary of the complexity:

- Applying the Master theorem (2.95) to the sort part yields a complexity in $\mathcal{O}(n \log n)$
- The complexity of the merge part is in $\mathcal{O}(n)$

Conclusion: the overall complexity of Sort and Count is $\mathcal{O}(n \log n)$



So far we have presented a few problems as well as some algorithms to solve them. We have also introduced the basics to study the complexity of those algorithms.

We now formalize these ideas through the following definitions.

Definitions

- ① A *computational problem* is a question or set of questions that a computer might be able to solve.
- ② The study of the solutions to computational problems composes the field of *Algorithms*.
- ③ *Computational complexity* attempts to classify algorithms depending on their speed or memory usage.

Taxonomy of computational problems

Decision problem: computational problem admitting exactly one of the two answers, “Yes” or “No”, to the question.

Example.

Let n be an integer, is n prime?

Search problem: computational problem where the answer is an arbitrary string.

Example.

Let n be an integer, find all the primes less than n .

Counting problem: computational problem where the answer is the number of solutions to a corresponding search problem.

Example.

Let n be an integer, count the number of primes less than n .

Taxonomy of computational problems

Optimization problem: computational problem where the answer is the best solution, with respect to some parameters, to a corresponding search problem.

Example.

Let n be a non-prime integer, find the largest prime factor of n .

Function problem: computational problem admitting exactly one answer for every input. The answer is more complex than in the case of a decision problem.

Example.

Given a list of cities and the distance between each pair, find the shortest route passing through all the cities exactly once and returning to the first visited city.

Relation between computational problems

An obvious observation is that counting and optimization problems are closely related to search problems.

As a less obvious observation one can notice that any optimization problem can be transformed into a decision problem.

Relation between computational problems

An obvious observation is that counting and optimization problems are closely related to search problems.

As a less obvious observation one can notice that any optimization problem can be transformed into a decision problem.

Example.

Optimization problem: let G be a graph and v_1, v_2 be two vertices in G . Find the shortest path between v_1 and v_2 .

A corresponding decision problem: let G be a graph and v_1, v_2 be two vertices in G . Is there a path from v_1 to v_2 that goes through less than 5 edges?

Relation between computational problems

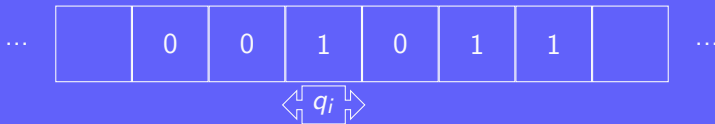
Similarly any function problem can be transformed into a decision problem. If the decision problem can be solved then so is its corresponding function problem. However the computational cost is not always preserved; for instance a function problem might be solved by an exponential time algorithm while its corresponding decision problem can be solved in polynomial time.

Conversely decision problems can be converted into function problems by computing the characteristic function of the set associated with the decision problem.

Decision problems play a central role in computability and complexity theories. In order to study them in more details we first need to setup a more precise computational model.

...an unlimited memory capacity obtained in the form of an infinite tape marked out into squares, on each of which a symbol could be printed. At any moment there is one symbol in the machine; it is called the scanned symbol. The machine can alter the scanned symbol, and its behavior is in part determined by that symbol, but the symbols on the tape elsewhere do not affect the behavior of the machine. However, the tape can be moved back and forth through the machine, this being one of the elementary operations of the machine. Any symbol on the tape may therefore eventually have an innings.

Turing – *Intelligent Machinery*



A Turing machine can be used to model modern computers:

- It is composed of an infinite tape
- The tape is divided into cells
- Each cell contains a symbol taken from an alphabet or a blank
- The tape is read or written, cell by cell, by a head
- At any time the device is in a state q_i with $i \in \mathbb{N}$

Let Σ be the alphabet of symbols and Q be the set of the states. Any elementary operation is determined by the symbol read and the state $q_i \in Q$.

Three basic cases can occur:

- The symbol read is replaced by another symbol from $\Sigma \cup \{b\}$, where b represents a blank
- The head moves to the left, right, or stands still
- The machine transitions from state q_i to state q_j

This can be summarized by the following map:

$$M : (\Sigma \cup \{b\}) \times Q \longrightarrow (\Sigma \cup \{b\}) \times \{-1, 0, 1\} \times Q$$

A deterministic computation is defined by (i) an initial state q_0 and (ii) a finite sequence of elementary operations.

At the end of the sequence the tape contains an integer x written over a number of cells, all the other cells being set to b (blank). The state of the tape, x , provides the result.

A deterministic computation is defined by (i) an initial state q_0 and (ii) a finite sequence of elementary operations.

At the end of the sequence the tape contains an integer x written over a number of cells, all the other cells being set to b (blank). The state of the tape, x , provides the result.

Definition

A function $f : \Sigma^* \longrightarrow \Sigma^*$ is said to be *Turing computable* if there exists a Turing machine M which returns $f(x)$ for any input x .

Remark.

Church proved that any computation, in a physical sense, can be performed on a Turing machine. This is however not the case anymore for quantum computers which cannot be modeled by a Turing machine. Its quantum counterpart is called *Quantum Turing machine*.

The RAM consists of:

- **A control unit:** containing a program and a program register pointing to the instruction to be executed
- **An arithmetic unit:** executes all the arithmetic operations
- **A memory:** divided into cells, each containing an integer
- **An input unit:** an input tape divided into cells and a head which reads the input
- **An output unit:** an output tape divided into cells and a head which writes the output

A program in the RAM model

The initial configuration of a RAM:

- All the cells are set to 0, besides the first n input cells
- The program register contains 1
- The first n cells contain the value from the n input cells

A computation consists in performing a sequence of configuration based on the program.

A program is composed of operands stored in the memory and onto which instructions can be run (e.g. LOAD, STORE, +, -, \times , /, READ, WRITE, JUMP, JZERO, JGE, HALT, ACCEPT, REJECT).

The RAM represents a better model for modern computer than a Turing machine. However those two models are equivalent.

Theorem

- For every Turing machine M , there exists a program P for RAM that simulates M .
- For a program P of RAM, there exists a Turing machine with five tapes such that P and M behave the same.

As we have briefly introduced the most common computational models we can now formalize the idea of complexity, and then investigate decision problems.

Definitions

Let Σ be an alphabet and M be a Turing machine.

- ① Given an input $x \in \Sigma^*$ for M , the number of operations $t_M(x)$ necessary to perform the computation is called the *length of the computation*.
- ② The function

$$\begin{aligned} T_M : \mathbb{N} &\longrightarrow \mathbb{N} \\ x &\longmapsto \max_{|x|} t_M(x) \end{aligned}$$

defines the *time complexity* for M .

Definitions

Let Σ be an alphabet and M be a Turing machine.

- 1 Given an input $x \in \Sigma^*$ for M , the number of operations $t_M(x)$ necessary to perform the computation is called the *length of the computation*.
- 2 The function

$$\begin{aligned} T_M : \mathbb{N} &\longrightarrow \mathbb{N} \\ x &\longmapsto \max_{|x|} t_M(x) \end{aligned}$$

defines the *time complexity* for M .

Remark.

Since the maximum of the length of the computation is considered this definition corresponds to the worst case complexity (2.68) with $|x|$, the size of the input x equal to n .

Definitions

- ① Let f be a Turing computable function, and M be that Turing machine. If there exists a polynomial P such that for any input x , $T_M(x) \leq P(x)$, then M is called a *deterministic polynomial algorithm*.
- ② Let $f : \Sigma^* \rightarrow \{0, 1\}$ be a Turing computable function and L be the preimage of $\{1\}$ under f , i.e. $f^{-1}(1)$. Then L is called a *language* and f defines a *decision problem* P . One says that the Turing machine *computes* f , *solves* P , or *decides* L .

Definitions

- ① The set of the decision problems which can be solved by a deterministic polynomial algorithm defines the *class* \mathcal{P} .
- ② A decision problem Π is computable by a non deterministic polynomial algorithm if and only if there exists a Turing machine M and a polynomial P such that
 - i $x \in L(\Pi)$ if and only if there exists $y \in \Sigma^*$, such that M computes 1 when it has input x in the cells 1 to $|x|$ and y in the cells -1 to $-|y|$.
 - ii For all x in $L(\Pi)$ there exists y such that M computes 1 in time less than $P(|x|)$.

The set of all such decision problems defines the class \mathcal{NP} .

Informal view:

- Class \mathcal{P} : set of decision problems solvable in polynomial time
- Class \mathcal{NP} : set of decision problems which (i) can be solved in polynomial time, assuming a *certificate* $y \in \Sigma^*$, is known and (ii) have True as answer.

The certificate is often taken to be the solution and it means that the solution can be verified in polynomial time.

- Class $\text{co-}\mathcal{NP}$: set of decision problems which (i) can be solved in polynomial time, assuming a *certificate* $y \in \Sigma^*$, is known and (ii) have False as answer.

The certificate is often taken to be the solution and it means that the solution can be verified in polynomial time.

- Any problem in \mathcal{P} is also in \mathcal{NP} and $\text{co-}\mathcal{NP}$.

Definitions

- ① Let P_1 and P_2 be two decision problems. One says that P_1 can be *reduced in polynomial time* to P_2 and writes $P_1 \times P_2$, if there exists a function f , computable in polynomial time, such that $x \in L(P_1)$ if and only if $f(x) \in L(P_2)$.
- ② A problem Π is \mathcal{NP} -hard if and only if for all P in \mathcal{NP} , P can be reduced in polynomial time to Π .
- ③ A problem Π is \mathcal{NP} -complete if and only if Π is (i) in \mathcal{NP} and (ii) \mathcal{NP} -hard.

Informal view:

- Saying that P_1 can be reduced in polynomial time to P_2 means that P_1 is not any harder than P_2 . In fact any instance of P_1 can be transformed, in polynomial time, into an instance of P_2 .
- An \mathcal{NP} -hard problem Π is at least as hard as the hardest problem in \mathcal{NP} . Note that Π does not need to belong to \mathcal{NP} , it could for instance be a search problem or an optimization problem.
- An \mathcal{NP} -complete problem Π must belong to \mathcal{NP} and is as hard as the hardest problem in \mathcal{NP} . In other words if we can solve Π then we can solve any other problem in \mathcal{NP} , deriving the solution from the one of Π , in polynomial time.

Problem (Halting problem)

Given the description of a Turing machine M as well as its initial input, decide whether M will halt or run forever.

The problem consists in constructing a Turing machine which would be able to decide whether or not another Turing machine would complete its task. From a simple perspective, if the machine “quickly” completes its task then we know it. Otherwise we have no idea whether it will halt later or never stop.

Turing proved that the halting problem is *undecidable*, i.e. it is impossible to construct an algorithm which always returns the right “yes” or “no” answer. Therefore it does not belong to \mathcal{NP} .

Problem (Boolean Satisfiability Problem (SAT))

Given a boolean function $(x_1, \dots, x_n) \mapsto f(x_1, \dots, x_n)$, where the x_i , $0 \leq i \leq n$ are in $\{0, 1\}$, and f is an expression constructed from the x_i and the boolean symbols \neg , \vee , and \wedge , is there a value for the n -tuple $\langle x_1, \dots, x_n \rangle$ such that $f(x_1, \dots, x_n)$ is True?

Example.

The expression $x_1 \wedge \neg x_2$ is satisfiable since taking $x_1 = \text{True}$ and $x_2 = \text{False}$ yields $\text{True} \wedge \text{True} = \text{True}$.

On the other hand the expression $x_1 \wedge \neg x_1$ is not satisfiable since it will always evaluate as False, independently of the choice of x_1 .

Theorem (Cook)

SAT is \mathcal{NP} -complete.

Sketch of proof.

First, SAT is in \mathcal{NP} as any instance can be verified in polynomial time on a Turing machine.

Second, any \mathcal{NP} problem P can be verified in polynomial time on a Turing machine M . Therefore for each input to M , it is possible to construct a boolean formula checking (i) every step of the computation, (ii) if M halts, and (iii) M returns “Yes”. Such a boolean expression will be satisfied if and only if those three conditions are met, that is if P is solved.

Finally as it can be proven that the boolean expression can be constructed in polynomial time, SAT is \mathcal{NP} -complete. \square

Theorem

The Halting problem is \mathcal{NP} -hard.

Sketch of proof.

The basic idea consists in reducing SAT (2.117) to the Halting problem (2.116).

It suffices to transform SAT into a Turing machine that tries all possible assignments of truth values. If a solution is found, then halt and otherwise start an infinite loop.

Since SAT is \mathcal{NP} -complete it means that any problem P in \mathcal{NP} can be reduced in polynomial time to the Halting problem. \square

Problem (True Quantified Boolean Formula (TQBF))

A quantified boolean formula is a formula that can be written in the form

$$Q_1x_1Q_2x_2\cdots Q_nx_n\phi(x_1, x_2, \dots, x_n),$$

where each of the Q_i is one of the quantifier \exists or \forall .

Calling L the language composed of the True quantified boolean formulae, decide L .

Example.

The quantified boolean formula $\forall x_1\exists x_2\forall x_3 ((x_1 \vee \neg x_2) \wedge (\neg x_1 \wedge x_3))$ is False.

Problem (True Quantified Boolean Formula (TQBF))

A quantified boolean formula is a formula that can be written in the form

$$Q_1x_1 Q_2x_2 \cdots Q_nx_n \phi(x_1, x_2, \cdots, x_n),$$

where each of the Q_i is one of the quantifier \exists or \forall .

Calling L the language composed of the True quantified boolean formulae, decide L .

Example.

The quantified boolean formula $\forall x_1 \exists x_2 \forall x_3 ((x_1 \vee \neg x_2) \wedge (\neg x_1 \wedge x_3))$ is False.

Observe that x_3 must be True in order for the formula to evaluate as True. However x_3 is a universally quantifiable variable.

Theorem

TQBF can be solved in exponential time and polynomial space.

Sketch of proof.

The idea consists in exhibiting a polynomial-space algorithm that decides whether any given Quantified Boolean Formula (QBF) is true.

Consider a general input

$$\Psi = Q_1 x_1 Q_2 x_2 \cdots Q_n x_n \phi(x_1, x_2, \cdots, x_n),$$

where ϕ has m clauses. We want to use the same argument as in example 2.120, i.e. if Q_i is \forall then both formulae, with x_i replaced by 0 and 1, must evaluate as True; if Q_i is \exists then only one of them is required to be true.

True quantified Boolean formula

Sketch of proof (continued).

This is achieved by “pilling off” the quantifiers one by one. Starting with Q_1 , one evaluates both

$$\psi_{1_0} = Q_2 x_2 \cdots Q_n x_n \phi(0, x_2, \cdots, x_n) \text{ and}$$

$$\psi_{1_1} = Q_2 x_2 \cdots Q_n x_n \phi(1, x_2, \cdots, x_n).$$

Then compute the logical value of $\psi_{1_0} \vee \psi_{1_1}$ or $\psi_{1_0} \wedge \psi_{1_1}$, depending whether Q_1 is \exists or \forall .

Following this strategy all the quantifiers are recursively removed and it then suffices to solve a SAT problem.

The computational cost of such a strategy is very high since at each recursive call the problem is transformed into two new linearly smaller sub-problems, resulting in a complexity of $\mathcal{O}(2^n)$.

Sketch of proof (continued).

We now consider the space necessary to solve TQBF.

As Ψ_i is computed either from $\Psi_{(i+1)_0} \wedge \Psi_{(i+1)_1}$ or $\Psi_{(i+1)_0} \vee \Psi_{(i+1)_1}$, the two instances of $\Psi_{(i+1)}$ can be computed sequentially. Therefore they can use the same memory space.

However at each level of recursion the indices must be saved. This incurs an $\mathcal{O}(\log n)$ overhead, and as a result if computing Ψ_{i+1} requires space S_{i+1} then evaluating Ψ_i requires

$$S_i = S_{i+1} + \mathcal{O}(\log n).$$

Hence Ψ can be reached in $\mathcal{O}(n \log n)$ space, that is in polynomial space.



Definition

- 1 Let SPACE denote the set of all the decision problems which can be solved by a Turing machine in $\mathcal{O}(s(n))$ space for some function s of the input size n . Then PSPACE is defined as

$$\text{PSPACE} = \bigcup_k \text{SPACE}(n^k).$$

- 2 A problem Π is *PSPACE-complete* if and only if (i) Π is in PSPACE and (ii) for all P in PSPACE , P can be reduced in polynomial space to Π .

Definition

- ① Let SPACE denote the set of all the decision problems which can be solved by a Turing machine in $\mathcal{O}(s(n))$ space for some function s of the input size n . Then PSPACE is defined as

$$\text{PSPACE} = \bigcup_k \text{SPACE}(n^k).$$

- ② A problem Π is *PSPACE-complete* if and only if (i) Π is in PSPACE and (ii) for all P in PSPACE , P can be reduced in polynomial space to Π .

Theorem

TQBF is PSPACE-complete and in particular it is \mathcal{NP} -hard.

Hilbert's tenth problem

Given a Diophantine equation with any number of unknown quantities and with rational integral numerical coefficients: To devise a process according to which it can be determined in a finite number of operations whether the equation is solvable in rational integers.

Hilbert – *List of Hilbert's problem list of 1900*

Hilbert's tenth problem

Given a Diophantine equation with any number of unknown quantities and with rational integral numerical coefficients: To devise a process according to which it can be determined in a finite number of operations whether the equation is solvable in rational integers.

Hilbert – *List of Hilbert's problem list of 1900*

Problem (Hilbert's tenth problem)

Given a Diophantine equation with any number of unknown quantities and with rational integral numerical coefficients, decide whether the equation is solvable in rational integers.

After a 21 years long quest Matiyasevich proved that Hilbert's tenth problem was undecidable, that is there is no algorithm which can determine whether any given Diophantine equation is solvable in rational integers.

A side result however states that one can exhibit a Diophantine equation with no solution, but such that it is impossible to prove it.

The idea is to number all the Diophantine equations, E_1, \dots, E_n, \dots following some set of axioms such as Peano arithmetic or Zermelo-Fraenkel set theory. The number of axioms and logical operators being finite one can construct all the possible proof P_1, \dots, P_n, \dots .

At stage 1, consider E_1 and check if:

- There exists a solution for the positive integers less than 1
- P_1 is a proof that E_1 has no solution

At stage k , consider E_1, \dots, E_k and check if:

- There exists a solution for the positive integers less than k
- P_1, \dots, P_k proves one of E_1, \dots, E_k

Then there is at least one Diophantine equation which has neither a solution nor a proof.

This conclusion is drawn from *Gödel incompleteness Theorem* which states that for any axiomatic system built over the Peano arithmetic there is an undecidable problem.

In order to relate *Gödel incompleteness Theorem* to Turing machine we informally define what it means for a set of axioms to be *complete* and *consistent*.

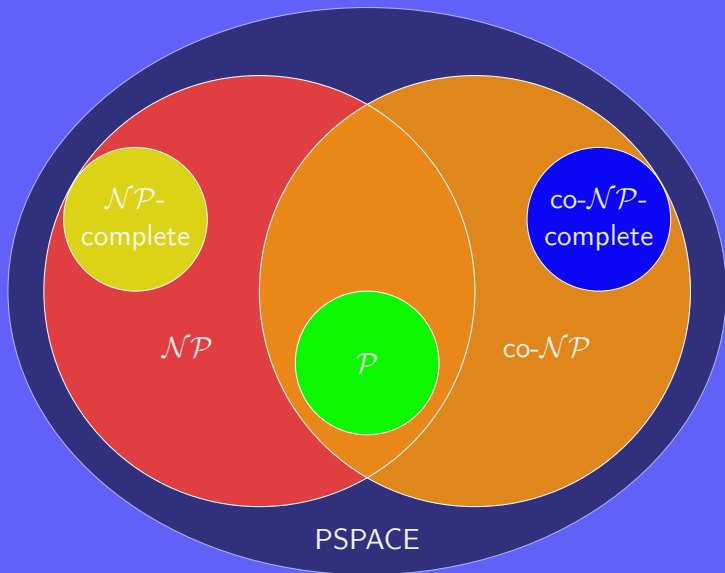
A set of axioms where any statement, built from it, or its negation can be proven in this set of axioms is said to be *complete*. It is *consistent* if it is impossible to prove both a statement and its negation in the set of axioms.

Assume we have a complete and consistent system F , powerful enough to reason about Turing machines. Taking a Turing machine M we can determine whether M halts by enumerating all the possible proofs in F , until one proof states that M halts or runs forever.

The completeness ensures of the correctness of the result while the consistency confirms the truthfulness of the conclusion.

We have proven that given F we can decide the Halting problem (2.116), which is known to be undecidable. ⚡

Therefore such a complete and consistent system F cannot exist.



Results known to be true:

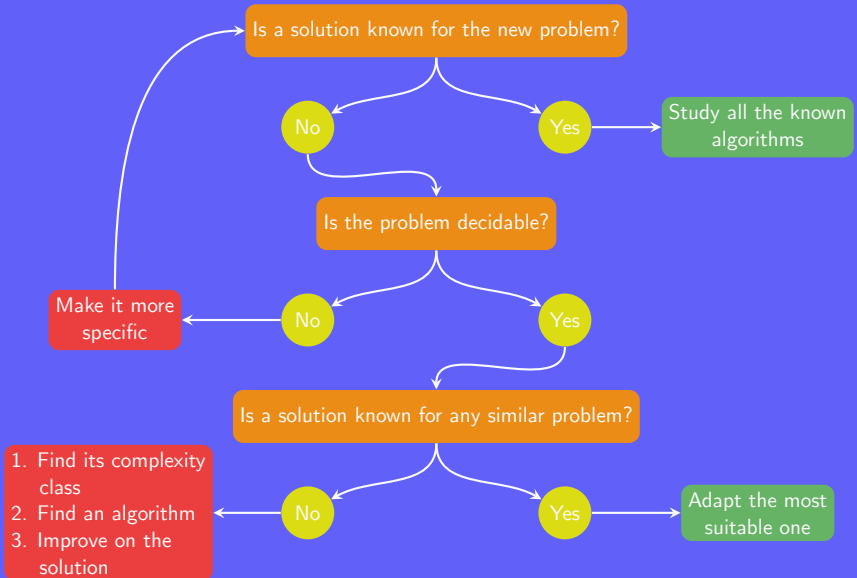
- $\mathcal{P} \subset \mathcal{NP}$
- $\mathcal{NP} \subset \text{PSPACE}$
- $\text{PSPACE} = \text{NPSPACE}$ (Savitch's theorem)

Results believed to be true:

- $\mathcal{P} \neq \mathcal{NP}$
- $\mathcal{NP} \neq \text{PSPACE}$

More results on computational complexity theory are available at the [Complexity Zoo](#).

Connection to algorithmic



We now provide a formal reduction proof as an example of how to proceed to evaluate the complexity class of a given problem. We start with the following definition.

Definition

A boolean formula is said to be in *Conjunctive Normal Form (CNF)* if it is written as the conjunction of disjunctive clauses.

As a reminder the truth table of the conjunction and disjunction of A and B is

A	B	$A \wedge B$	$A \vee B$
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	1

Problem (Satisfiability with 3 literals per clause (3-SAT))

Given a Boolean formula in CNF where each clause contains exactly three literals is there a truth assignment which satisfies all the clauses.

Theorem

3-SAT is \mathcal{NP} -complete.

Proof.

Clearly 3-SAT is in \mathcal{NP} , for it is a particular case of SAT.

To prove that 3-SAT is \mathcal{NP} -hard we will show that being able to solve it implies being able to solve SAT, which by Cook theorem (2.118) is known to be \mathcal{NP} -hard.

Proof (continued).

Assuming CNF, we want to transform any instance of SAT into an instance of 3-SAT. Therefore we need to consider the cases where SAT has clauses with (i) one, (ii) two or (iii) more than three literals. Note that the case where there are exactly three clauses is already 3-SAT so no work is necessary.

(ii) The most simple case is when there are two literals, organised in a unique disjunctive clause $x_1 \vee x_2$ denoted (x_1, x_2) . It then suffices to consider the pair of (x_1, x_2, u) and $(x_1, x_2, \neg u)$, where u is a newly added literal.

(i) Similarly in the case of a single literal x we convert it into a pair of literals (x, u_1) and $(x, \neg u_1)$, for a new literal u_1 . It then suffices to apply the same strategy as above and create the four literals

$$(x, u_1, u_2), (x, u_1, \neg u_2), (x, \neg u_1, u_2), (x, \neg u_1, \neg u_2).$$

Proof (continued).

The general idea behind (i) and (ii) is to add new literals which do not impact the satisfiability of the problem.

(iii) The case of a clause with more than three literals (x_1, \dots, x_n) can be treated by adding some new literals such as to create a chain of clauses where each of them has exactly three literals.

$$(x_1, x_2, u_3), (\neg u_3, x_3, u_4), (\neg u_4, x_4, u_5), \dots, (\neg u_{n-1}, x_{n-1}, x_n).$$

The question is now to know if this change alters the satisfiability of the original clause.

Assume the original clause evaluates as True. Then at least one literal x_i must be True. Setting u_j to True for all $j \leq i$ and False for $j > i$ preserves the satisfiability of the clause.

Proof (continued).

If now the original clause evaluates as `False` we need to ensure that the altered version cannot be satisfied. If it is originally `False` then all the x_i must be `False`.

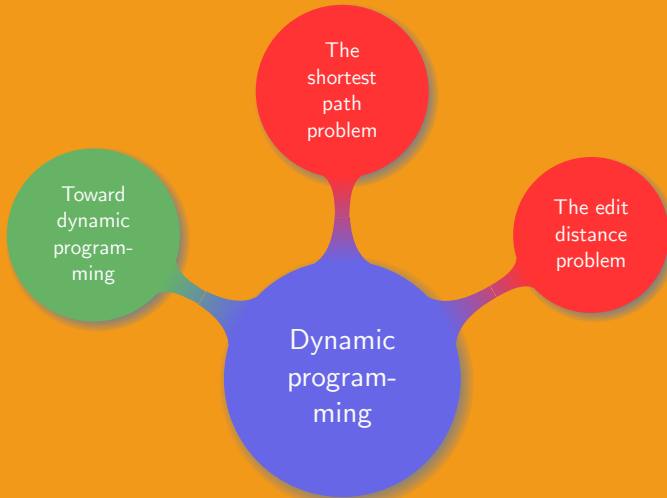
Therefore in the altered version the last clause can only be `True` if u_{n-1} is `False`. In turn this implies that u_{n-2} must also be `False`, and by extension all the u_i for $3 \leq i \leq n-2$. However in that case the first clause fails as u_3 is `False` and this means the conjunction of all the clauses evaluates as `False`.

In order to finalize the proof we need to ensure that the proposed transformations are applicable in polynomial time. This is clearly the case since the number of additional literals is polynomial in the number of original literals.

Hence $SAT \times 3\text{-SAT}$ and 3-SAT is \mathcal{NP} -complete.



3. Dynamic programming



Fibonacci posed the following problem in his book *Liber Abbaci* (Book of Calculations) in 1202:

A certain man puts a pair of rabbits in a place surrounded on all sides by a wall. How many pairs of rabbits can be produced from that pair in a year if it is supposed that every month each pair begets a new pair, which from the second month on becomes productive?

The solution leads to the sequence of Fibonacci numbers:

- At the beginning of month 1, there is $F_0 = 1$ pair, which is not productive.
- At the beginning of month 2, there is still $F_1 = 1$ pair, which is now productive.
- At the beginning of month 3, there are now $F_2 = 2$ pairs, of which one is productive.

- At the beginning of month 4, there are now $F_3 = 3$ pairs and the pair born in month 3 becomes productive.
- As the beginning of month 5, the number of pairs is equal to those of month 4, plus all those that were productive in month 4. These are all pairs that existed in month 2, since all of those will be productive in month 3. Hence, $F_4 = 3 + 2 = 5$.
- In general, at the beginning of every month the number of pairs of rabbits is equal to the number of pairs of the previous month, plus the number of pairs of two months ago, which have since become productive.

Hence,

$$F_n = F_{n-1} + F_{n-2}, \quad n \geq 2, \quad F_0 = 1, \quad F_1 = 1.$$

Algorithm.

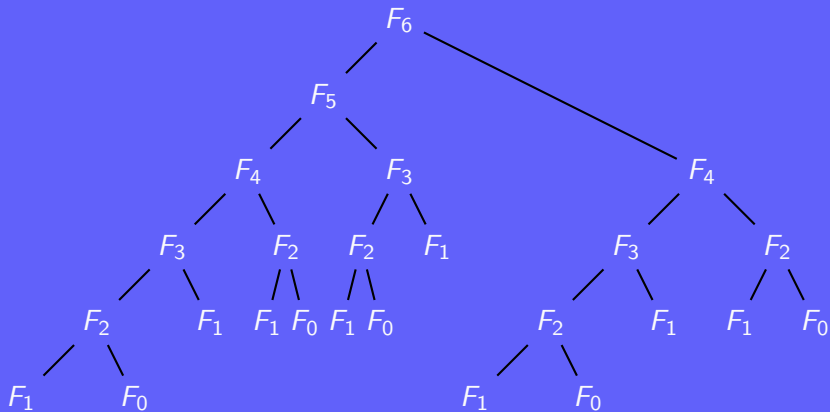
Input : An integer n

Output: F_n

```
1 Function Fib_vn( $n$ ):  
2   if  $n = 0$  then return 0;  
3   if  $n = 1$  then return 1;  
4   return Fib_vn( $n - 1$ ) + Fib_vn( $n - 2$ )  
5 end
```

Since $F_{n+1}/F_n \approx \Phi = \frac{1+\sqrt{5}}{2} > 1.6$, it means that $F_n > 1.6^n$.

Noting that each recursive call has to reach 0 or 1 to be completed it is clear that the complexity of this algorithm is exponential.



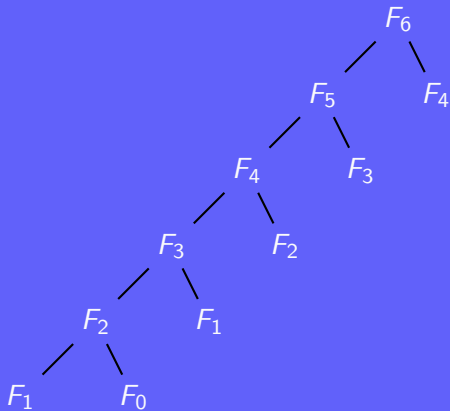
Algorithm.

Input : An integer n

Output: F_n

```
1 Function Fib_n( $n$ ):  
2   if  $F_n = -1$  then  $F_n \leftarrow \text{Fib\_n}(n-1) + \text{Fib\_n}(n-2)$ ;  
3   return  $F_n$   
4 end  
5 Function Fib_nm( $n$ ):  
6    $F_0 \leftarrow 0$ ;  
7    $F_1 \leftarrow 1$ ;  
8   for  $i \leftarrow 2$  to  $n$  do  $F_i \leftarrow -1$ ;  
9   return Fib_n( $n$ )  
10 end
```

Since each value is computed exactly once and saved, both the time and space complexities are linear in n .



Algorithm.

Input : An integer n

Output: F_n

```
1 Function Fib( $n$ ):  
2    $F_{old2} \leftarrow 0; F_{old1} \leftarrow 1;$   
3   if  $n = 0$  then return 0;  
4   for  $i \leftarrow 2$  to  $n$  do  
5      $F \leftarrow F_{old1} + F_{old2}; F_{old2} \leftarrow F_{old1}; F_{old1} \leftarrow F;$   
6   end for  
7   return  $F_{old1} + F_{old2}$   
8 end
```

With this simple strategy the time is linear in n , while the storage is constant.

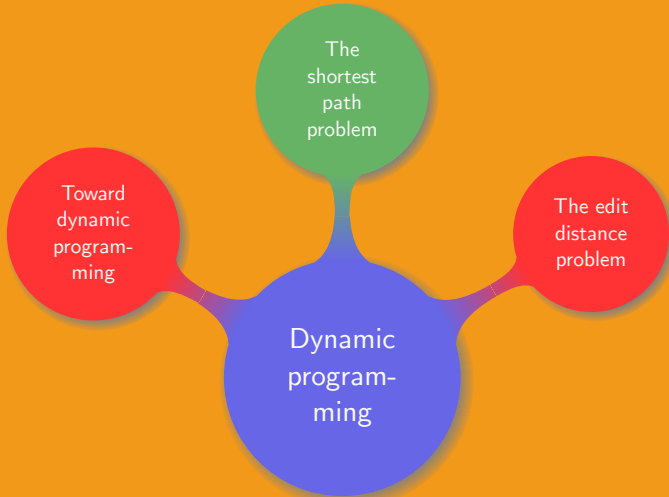
Simple idea behind dynamic programming:

- Break a complex problem into simpler subproblems
- Store the result of the overlapping subproblems
- Combine the solutions of the subproblems to solve the initial problem

Simple idea behind dynamic programming:

- Break a complex problem into simpler subproblems
- Store the result of the overlapping subproblems
- Combine the solutions of the subproblems to solve the initial problem

Dynamic programming is simply a *space-time trade-off*.



Problem (Shortest paths in weighted graphs)

Given a connected, simple, weighted graph, and two vertices s and t , find the shortest path that joins s to t .

Two main cases for the graph:

- It only has edges with positive labels
- It has edges with positive and negative labels

We now recall Dijkstra's algorithm to solve the first case and then introduce the Bellman-Ford algorithm which takes advantage of dynamic programming in order to treat the second one.

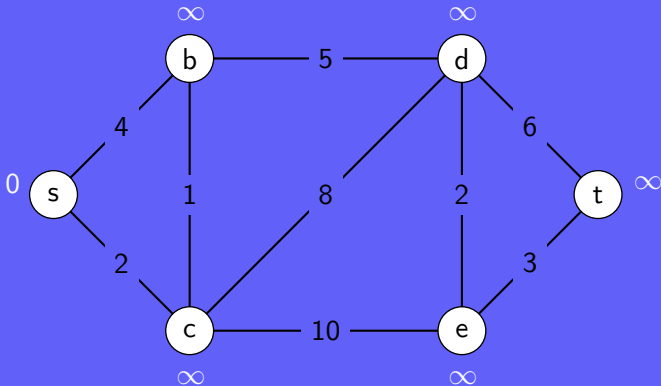
Algorithm. (*Dijkstra*)

Input : A graph $G = \langle V, E \rangle$ with positive edges, two vertices s and t

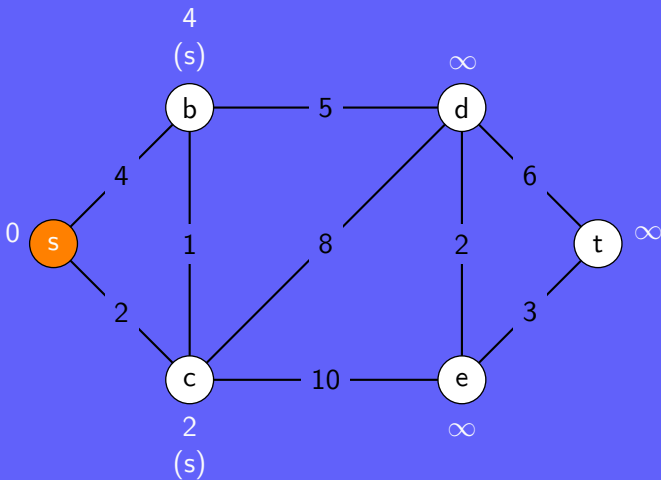
Output : The shortest path between s and t

```
1  $s.dist \leftarrow 0$ ;  $s.prev \leftarrow \text{NULL}$ ;  $S \leftarrow \emptyset$ ;  
2 foreach vertex  $v \in G.V$  do  
3   |   if  $v \neq s$  then  $v.dist \leftarrow \infty$ ;  $v.prev \leftarrow \text{NULL}$ ;  
4   |   add  $v$  to  $S$ ;  
5 end foreach  
6  $v \leftarrow s$ ;  
7 repeat  
8   |   foreach neighbor  $u$  of  $v$  do  
9   |   |    $tmp \leftarrow v.dist + \text{weight}(v, u)$ ;  
10  |   |   if  $tmp < u.dist$  then  $u.dist \leftarrow tmp$ ;  $u.prev \leftarrow v$ ;  
11  |   end foreach  
12  |   remove  $v$  from  $S$ ;  $v \leftarrow$  vertex with minimal distance in  $S$ ;  
13 until  $v = t$ ;  
14 return  $t, t.prev, \dots, s$ 
```

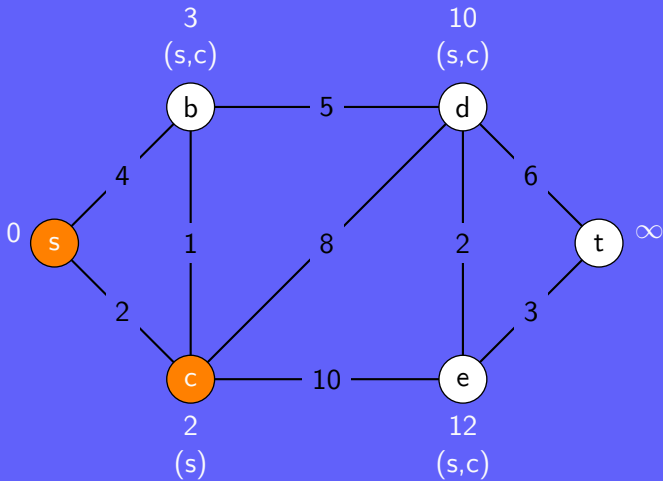
$$S_0 = \emptyset.$$



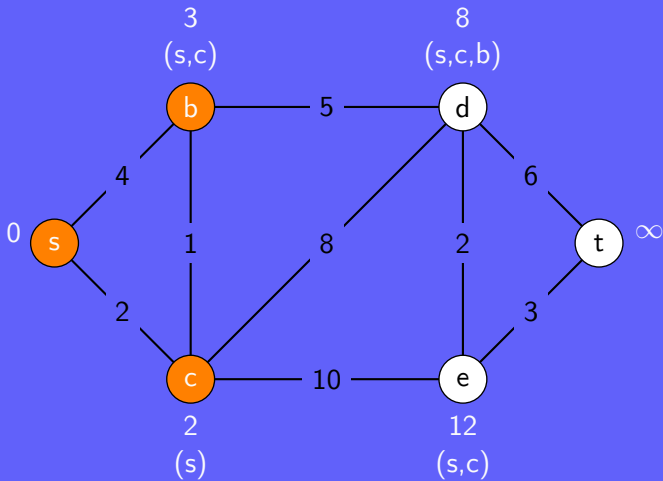
$$S_1 = \{s\}.$$



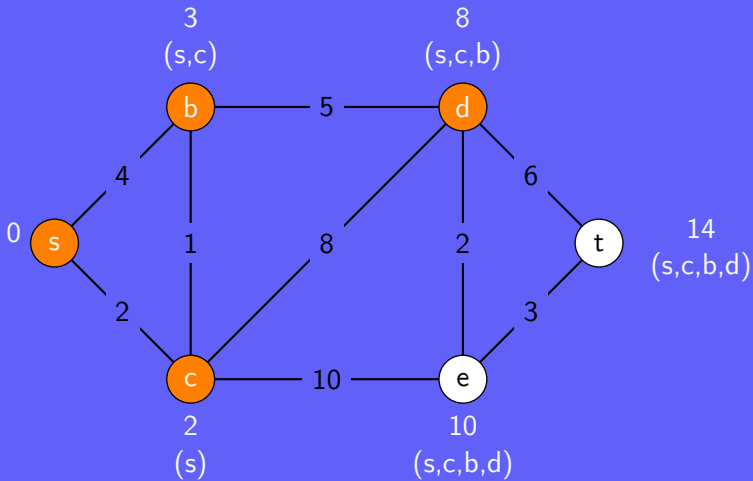
$$S_2 = \{s, c\}.$$



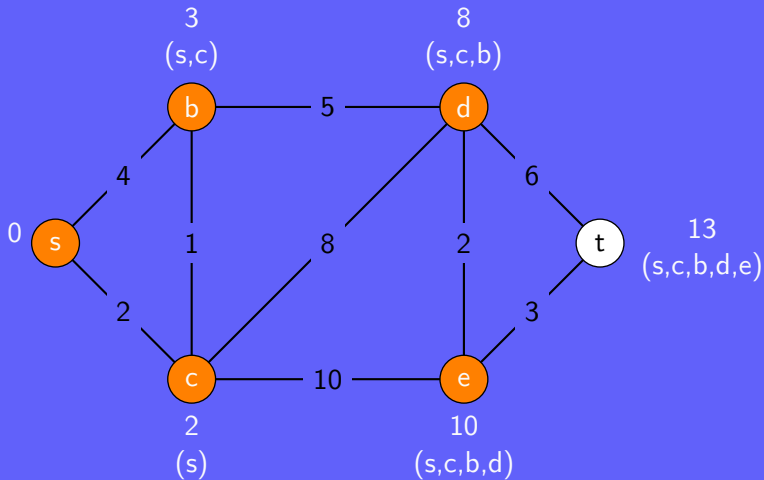
$$S_3 = \{s, b, c\}.$$



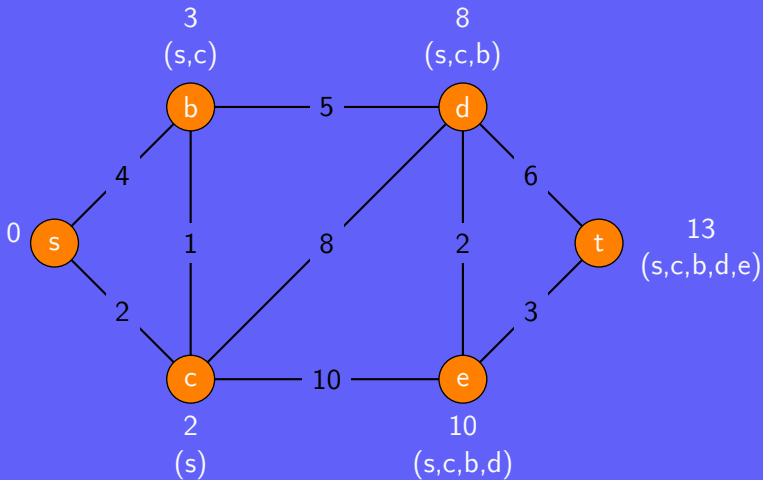
$$S_4 = \{s, b, c, d\}.$$



$$S_5 = \{s, b, c, d, e\}.$$



$$S_6 = \{s, b, c, d, e, t\}.$$



Shortest path with negative edges

Proposition

If G is a graph with no negative cycle, then there is a shortest simple path going from a source vertex s to a target vertex t .

Proof.

Assume no cycle of negative cost exists. If the path repeated a vertex then removing edges to break this cycle would result in a path of no greater cost and with fewer edges. Therefore a shortest simple path exists. \square

Given a graph with n nodes we can find a path of minimum cost which has length at most $n-1$. Let $opt(i, v)$ denote the minimum cost of path from a vertex v to a vertex t and featuring at most i edges. The goal is to express $opt(i, v)$ into smaller subproblems such as to determine $opt(n-1, s)$ using dynamic programming.

Shortest path with negative edges

Lemma

Let P be an optimal path $opt(i, v)$ in a graph G , and (v, w) be the first edge in P . Then

$$opt(i, v) = \min(opt(i-1, v), opt(i-1, w) + \text{weight}(v, w))$$

Proof.

The recursive formula is clear as soon as one notes that two cases can occur.

First if P has at most $i-1$ edges then $opt(i, v) = opt(i-1, v)$. On the other hand if P has i edges and the first one is (v, w) , then

$$opt(i, v) = \text{weight}(v, w) + opt(i-1, w).$$



Algorithm. (*Bellman-Ford*)

Input : A directed weighted graph $G = \langle V, E \rangle$ with no negative cycle, two vertices s and t

Output: The shortest path between s and t

```
1 foreach  $v \in G.V$  do
2   if  $v \neq s$  then  $v.dist \leftarrow \infty$ ;  $v.prev = \text{NULL}$ ;
3   else  $v.dist \leftarrow 0$ ;
4 end foreach
5 for  $i \leftarrow 1$  to  $|G.V| - 1$  do
6   foreach  $(u, v) \in G.E$  do
7      $tmp \leftarrow u.dist + \text{weight}(u, v)$ ;
8     if  $tmp < v.dist$  then  $v.dist \leftarrow tmp$ ;  $v.prev = u$ ;
9   end foreach
10 end for
11 return  $t, t.prev, \dots, s$ 
```

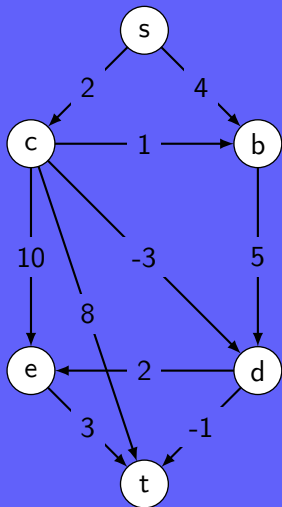
Theorem

Bellman-Ford algorithm is correct and runs in time $\mathcal{O}(mn)$ on a graph composed of n vertices and m edges.

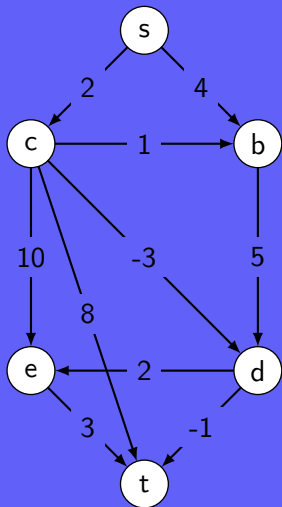
Proof.

The correctness of the algorithm follows from the recursive formula introduced in lemma 3.153.

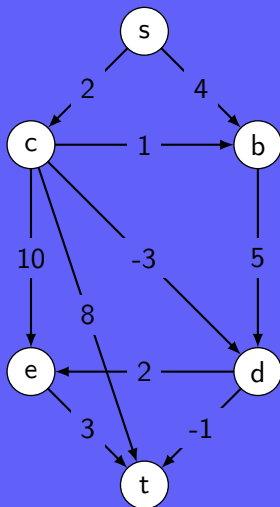
From lines 5 to 9 of Bellman-Ford algorithm (3.154) the loop on all the m edges is applied n times. □



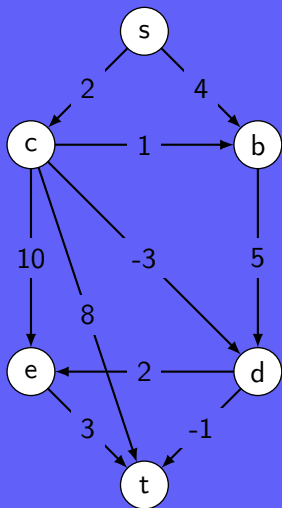
	0	1	2	3	4	5
<i>s</i>	0					
<i>b</i>	∞					
<i>c</i>	∞					
<i>d</i>	∞					
<i>e</i>	∞					
<i>t</i>	∞					



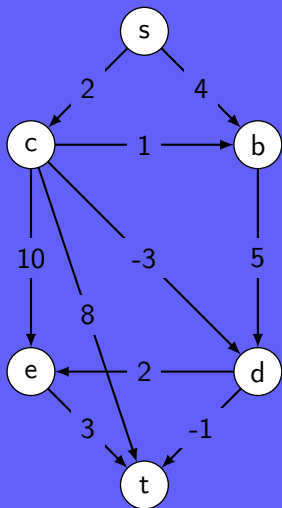
	0	1	2	3	4	5
<i>s</i>	0	0				
<i>b</i>	∞	4				
<i>c</i>	∞	2				
<i>d</i>	∞	∞				
<i>e</i>	∞	∞				
<i>t</i>	∞	∞				



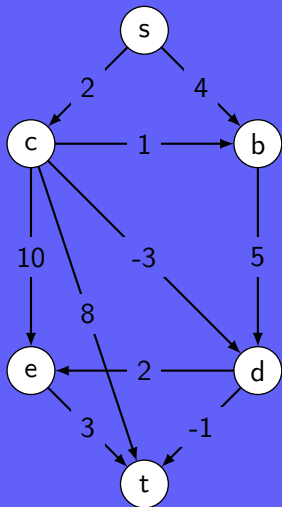
	0	1	2	3	4	5
<i>s</i>	0	0	0			
<i>b</i>	∞	4	3			
<i>c</i>	∞	2	2			
<i>d</i>	∞	∞	-1			
<i>e</i>	∞	∞	12			
<i>t</i>	∞	∞	10			



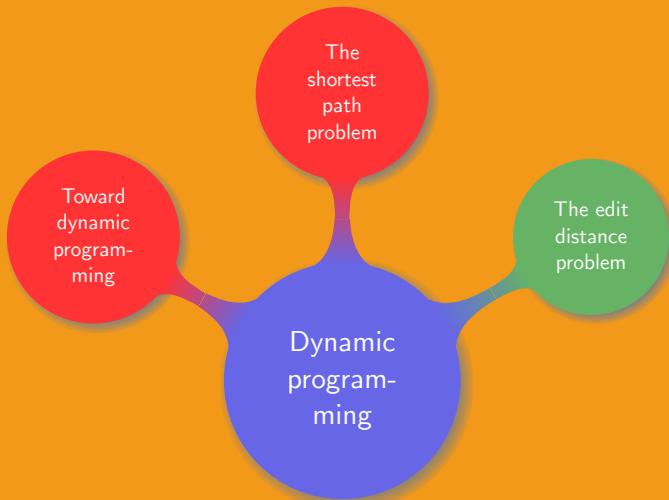
	0	1	2	3	4	5
<i>s</i>	0	0	0	0		
<i>b</i>	∞	4	3	3		
<i>c</i>	∞	2	2	2		
<i>d</i>	∞	∞	-1	-1		
<i>e</i>	∞	∞	12	1		
<i>t</i>	∞	∞	10	-2		



	0	1	2	3	4	5
<i>s</i>	0	0	0	0	0	
<i>b</i>	∞	4	3	3	3	
<i>c</i>	∞	2	2	2	2	
<i>d</i>	∞	∞	-1	-1	-1	
<i>e</i>	∞	∞	12	1	1	
<i>t</i>	∞	∞	10	-2	-2	



	0	1	2	3	4	5
<i>s</i>	0	0	0	0	0	0
<i>b</i>	∞	4	3	3	3	3
<i>c</i>	∞	2	2	2	2	2
<i>d</i>	∞	∞	-1	-1	-1	-1
<i>e</i>	∞	∞	12	1	1	1
<i>t</i>	∞	∞	10	-2	-2	-2



Basic idea: given two strings, do they match?

Comments on the problem:

- Simple to implement
- How to render misspellings?

Conclusion: useless in practice

Goal: measure how far two strings are from each others

Basic idea: given two strings, do they match?

Comments on the problem:

- Simple to implement
- How to render misspellings?

Conclusion: useless in practice

Goal: measure how far two strings are from each others

Typical applications:

- Spell checker
- Changes in language usage
- DNA sequencing
- Plagiarism

Problem (Edit distance)

The number of changes that need to be performed to convert a string into another one defines the *distance* between the two strings. The three types of alterations considered are:

- *Substitution*: a single character is replaced by a different one
- *Insertion*: a single character is added such as to decrease the distance between the two strings
- *Deletion*: a single character is deleted in order to match the other string

Given two strings and assigning distance one to each of these operations determine the *edit distance* between them.

Naive idea: search exact places where to add/delete characters

Alternative view:

- What information is needed to decide on what operation to perform?
- What can happen to the last character for each string?

Naive idea: search exact places where to add/delete characters

Alternative view:

- What information is needed to decide on what operation to perform?
- What can happen to the last character for each string?

If an optimal solution is known for all the characters but the last, then it becomes simple to find an overall best solution: check the three possibilities, add the cost of each of them to the previous minimal cost and select the best option.

Given a reference string S and a string T we want to determine their edit distance. At each step, i.e. letter, a decision can be taken upon the previous results:

- If $S_i = T_j$, then consider $dist_{i-1,j-1}$. Otherwise consider $dist_{i-1,j-1}$ and pay a cost 1 for the difference
- If $S_{i-1} = T_j$, then it could be that T has one more character than S . In that case consider $dist_{i-1,j}$ and pay a cost 1 for the insertion of a character in S
- If $S_i = T_{j-1}$, then it could be that T has one less character than S . In that case consider $dist_{i,j-1}$ and pay a cost 1 for the deletion of a character in S

As those three possibilities cover all the cases, taking their minimum yields the edit distance.

Limitations of the recursive approach

Description of the strategy:

- If either of the index is 0 then set *dist* to the other index
- Compute

$$dist_{i,j} = \min(dist_{i-1,j} + 1, dist_{i,j-1} + 1, dist_{i-1,j-1} + (\text{match}(S_i, T_j) ? 0 : 1))$$

How fast would be the algorithm?

Limitations of the recursive approach

Description of the strategy:

- If either of the index is 0 then set *dist* to the other index
- Compute

$$dist_{i,j} = \min(dist_{i-1,j} + 1, dist_{i,j-1} + 1, dist_{i-1,j-1} + (\text{match}(S_i, T_j) ? 0 : 1))$$

How fast would be the algorithm?

- At each position in the string, three branches are explored
- Only one branch reduces both indices
- Exponential time $\Omega(3^n)$

How to do better?

Limitations of the recursive approach

Description of the strategy:

- If either of the index is 0 then set $dist$ to the other index
- Compute

$$dist_{i,j} = \min(dist_{i-1,j} + 1, dist_{i,j-1} + 1, dist_{i-1,j-1} + (\text{match}(S_i, T_j) ? 0 : 1))$$

How fast would be the algorithm?

- At each position in the string, three branches are explored
- Only one branch reduces both indices
- Exponential time $\Omega(3^n)$

How to do better?

- What is the maximum number of pairs?
- The same pairs are recalled many times
- Use a lookup table to decrease the computational cost

Algorithm.

Input : Two strings S and T

Output: The edit distance between S and T

```
1 for  $i \leftarrow 0$  to  $|S|$  do  $dist_{i,0} \leftarrow i$ ;  
2 for  $i \leftarrow 1$  to  $|T|$  do  $dist_{0,i} \leftarrow i$ ;  
3 for  $i \leftarrow 1$  to  $|S|$  do  
4   for  $j \leftarrow 1$  to  $|T|$  do  
5      $tmp_0 \leftarrow dist_{i-1,j-1} + (\text{match}(S_i, T_j) ? 0 : 1)$ ;  
6      $tmp_1 \leftarrow dist_{i,j-1} + 1$ ;           /* skip a letter in  $S$  */  
7      $tmp_2 \leftarrow dist_{i-1,j} + 1$ ;       /* skip a letter in  $T$  */  
8      $dist_{i,j} \leftarrow \min(tmp_0, tmp_1, tmp_2)$ ;  
9   end for  
10 end for  
11 return  $dist_{i,j}$ 
```

		P	O	L	Y	N	O	M	I	A	L
	0	1	2	3	4	5	6	7	8	9	10
E	1	1	2	3	4	5	6	7	8	9	10
X	2	2	2	3	4	5	6	7	8	9	10
P	3	2	3	3	4	5	6	7	8	9	10
O	4	3	2	3	4	5	6	7	8	9	10
N	5	4	3	3	4	4	5	6	7	8	9
E	6	5	4	4	4	5	5	6	7	8	9
N	7	6	5	5	5	4	5	6	7	8	9
T	8	7	6	6	6	5	5	6	7	8	9
I	9	8	7	7	7	6	6	6	6	7	8
A	10	9	8	8	8	7	7	7	7	6	7
L	11	10	9	8	9	8	8	8	8	7	6

Theorem

Let S and T be two strings of length n and m , respectively. Algorithm 3.163 solves edit distance in time $\mathcal{O}(nm)$.

Proof.

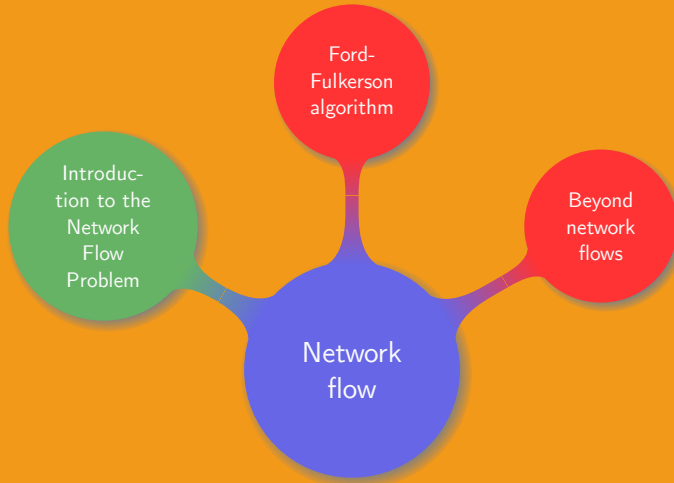
Algorithm 3.163 is simply a table look up version of the recursive algorithm described on slide 3.161. As all the cases are considered in the initial recursive approach they are also all covered in Algorithm 3.163.

Looking at the description of the algorithm it is clear that most of the work is performed in the nested for loops. This results in the creation of a $n \times m$ table. All the other operations are only reading from this lookup table. Therefore the complexity is $\mathcal{O}(nm)$. \square

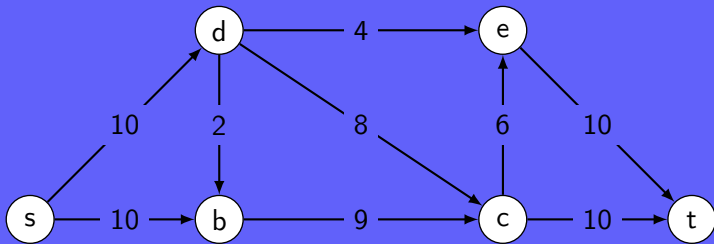
Dynamic programming:

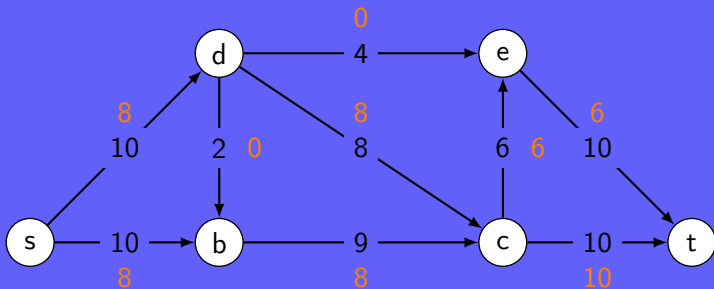
- Simple space-time trade-off
- Cover all the possibilities at the subproblems level
- Never recompute anything that is already known
- Efficient as long as the number of subproblems remains polynomial

4. Network flow

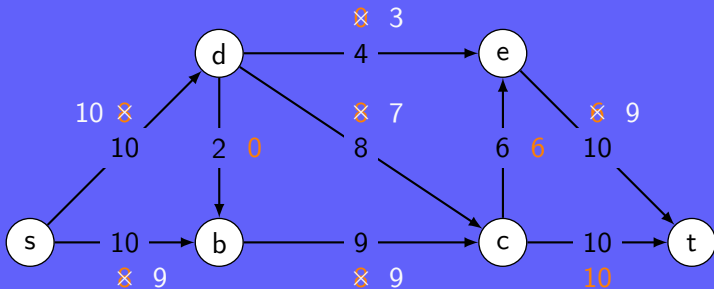








Total flow: 16



Total flow: 19

Let $G = \langle V, E \rangle$ be a directed graph. We consider each node of G as a switch and each edge as carrying some traffic. This is for instance the case in the context of a highway system: nodes are the interchanges (hubs) and the edges the highway itself. It could also be illustrated using a fluid network where the edges are the pipes and the nodes the junctures where the pipes are plugged together.

To each edge one associates a number called *capacity*, which represents how much traffic an edge can handle. In the maximum network flow problem the goal is to arrange the traffic such as optimizing the available capacity.

Before being able to solve this problem we need to formalize the idea of flow.

Definition

Let $G = \langle V, E \rangle$ be a weighted directed graph with a *source* node s and a *sink* node t ; G is called a *flow network*.

Given a function $c : E \rightarrow \mathbb{R}^+$, called *capacity*, a *flow* is a function $f : E \rightarrow \mathbb{R}^+$, satisfying the following properties.

- i *Capacity constraint*: the flow of an edge can never exceed its capacity, i.e. $\forall (u, v) \in E, f(u, v) \leq c(u, v)$.
- ii *Flow conservation*: at each node the entering flow equals the exiting flow, i.e.

$$\forall u \in V \setminus \{s, t\}, \sum_{v \in V} f(v, u) = \sum_{v \in V} f(u, v).$$

Definition

Let $G = \langle V, E \rangle$ be a weighted directed graph with a *source* node s and a *sink* node t ; G is called a *flow network*.

Given a function $c : E \rightarrow \mathbb{R}^+$, called *capacity*, a *flow* is a function $f : E \rightarrow \mathbb{R}^+$, satisfying the following properties.

- i *Capacity constraint*: the flow of an edge can never exceed its capacity, i.e. $\forall (u, v) \in E, f(u, v) \leq c(u, v)$.
- ii *Flow conservation*: at each node the entering flow equals the exiting flow, i.e.

$$\forall u \in V \setminus \{s, t\}, \sum_{v \in V} f(v, u) = \sum_{v \in V} f(u, v).$$

Remark.

Unless specified otherwise, we restrict our attention to the case of integer or at least rational capacities and flows.

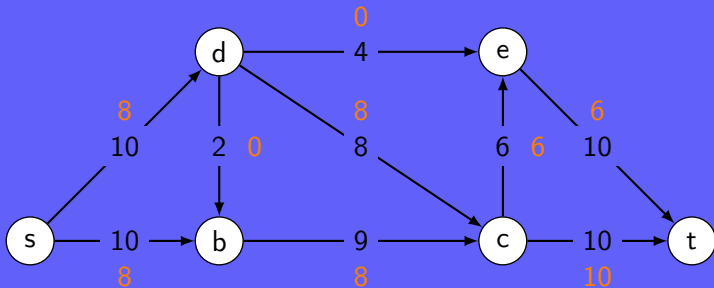
Problem (Maximum Network Flow)

Given a flow network arrange the flow such as to maximize the available capacity.

Remark.

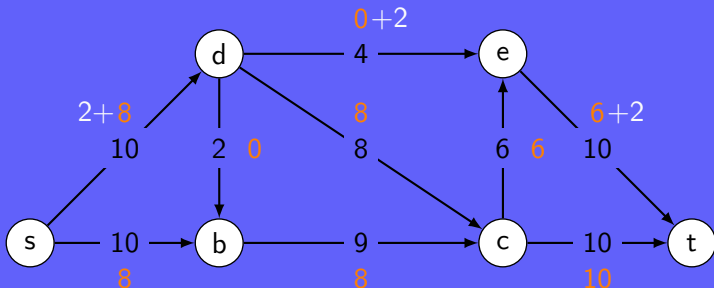
Assume a network flow is divided into two sets A and B of empty intersection and s is in A while t is in B . Intuitively any flow travelling from s to t must cross from A to B . This suggests that the capacity of the network flow can never exceed the capacity of the *cuts* A and B . The minimum capacity of any such division is called the *minimum cut*. We will prove that it is equal to the maximum flow value.

Back to the toy example



Improving the original path:

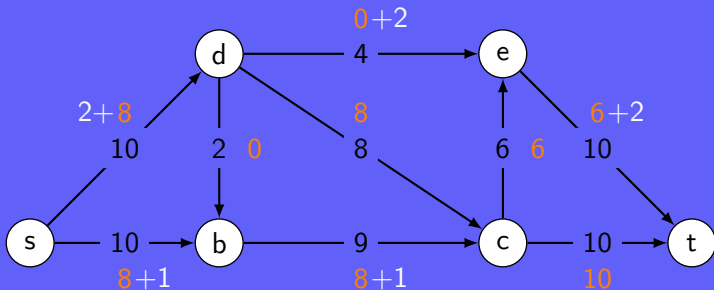
Back to the toy example



Improving the original path:

- Add $+2$ along $s \rightarrow d \rightarrow e \rightarrow t$; total flow: 18

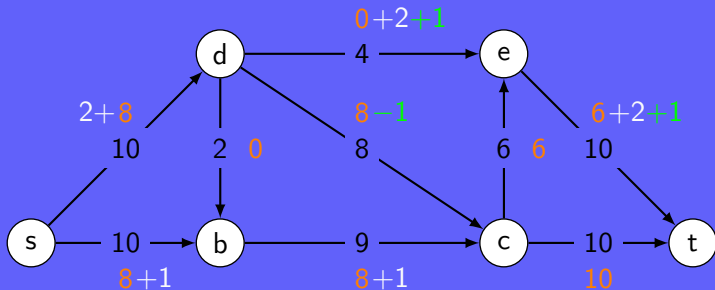
Back to the toy example



Improving the original path:

- Add +2 along $s \rightarrow d \rightarrow e \rightarrow t$; total flow: 18
- Add +1 along $s \rightarrow b \rightarrow c$: c is a bottleneck; total flow: 18

Back to the toy example



Improving the original path:

- Add $+2$ along $s \rightarrow d \rightarrow e \rightarrow t$; total flow: 18
- Add $+1$ along $s \rightarrow b \rightarrow c$: c is a bottleneck; total flow: 18
- Reallocate 1 from $d \rightarrow c \rightarrow t$ to $d \rightarrow e \rightarrow t$; total flow: 19

Finding the maximum flow can be achieved as follows:

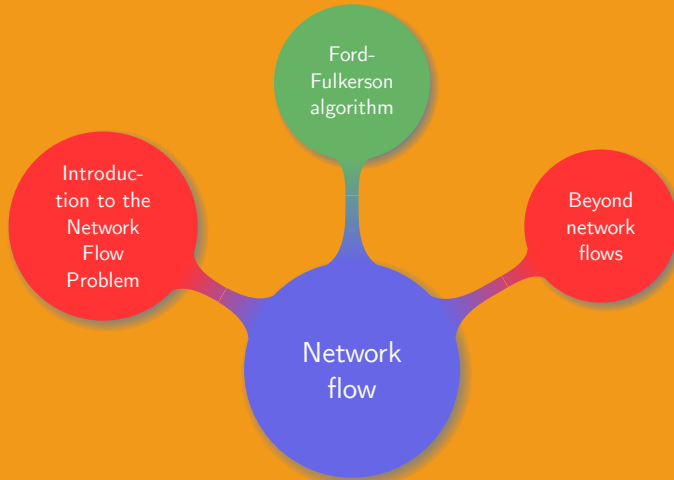
- ① Start with a null flow
- ② Find paths to increase the flow:
 - Path where the capacity has not been reached
 - Reallocate flow to a different path and increase the capacity of the current one
- ③ Stop when no more flow can be injected

Finding the maximum flow can be achieved as follows:

- ① Start with a null flow
- ② Find paths to increase the flow:
 - Path where the capacity has not been reached
 - Reallocate flow to a different path and increase the capacity of the current one
- ③ Stop when no more flow can be injected

Goals:

- Determine all the paths that allow an increase of the flow
- Make sure the search for such paths stops at some stage



Definition (Residual graph)

Let G be a graph and f be a flow network. The *residual graph* G_f of G with respect to f is the graph whose vertices are the vertices of G ; Regarding its edges:

- Each edge e of G with capacity $c(e)$ and flow $f(e) < c(e)$ is added to G_f with capacity $c(e) - f(e)$; it is a *forward edge*.
- For each edge $e = (u, v)$ of G with flow $f(e) > 0$, a *backward edge* $e' = (v, u)$ with capacity $c(e') = f(e)$ is added; it is a *backward edge*.

For a path P , the maximum amount by which the flow can be increased on each edge along P is called the *residual capacity*.

Definition (Residual graph)

Let G be a graph and f be a flow network. The *residual graph* G_f of G with respect to f is the graph whose vertices are the vertices of G ; Regarding its edges:

- Each edge e of G with capacity $c(e)$ and flow $f(e) < c(e)$ is added to G_f with capacity $c(e) - f(e)$; it is a *forward edge*.
- For each edge $e = (u, v)$ of G with flow $f(e) > 0$, a *backward edge* $e' = (v, u)$ with capacity $c(e') = f(e)$ is added; it is a *backward edge*.

For a path P , the maximum amount by which the flow can be increased on each edge along P is called the *residual capacity*.

Remark.

Given a graph G its residual graph G_f has at most twice the number of edges as G .

Algorithm. (*Augment*)

Input : a flow f and a simple path P

Output: a new flow

```
1 Function Augment( $f, P$ ):  
2    $b \leftarrow$  minimum residual capacity on  $P$  with respect to  $f$ ;  
3   foreach edge  $e \in P$  do  
4     if  $e$  is a forward edge then  $f(e) \leftarrow f(e) + b$ ;  
5     else  $f(e) \leftarrow f(e) - b$ ;  
6   end foreach  
7   return  $f$   
8 end
```

Lemma

The Augment algorithm (4.179) returns a flow in the original graph G .

Proof.

We must ensure that the output of the algorithm matches definition 4.173, that is satisfies the capacity constraint and the flow conservation properties.

Consider the edges from the residual graph G_f whose capacity differ from the ones of G . In the case of e being such a forward edge with capacity $c(e) - f(e)$, we have

$$0 \leq f(e) \leq f(e) + b \leq f(e) + (c(e) - f(e)) = c(e).$$

Proof (continued).

If e is a backward edge then its residual capacity is $f(e)$ and

$$c(e) \geq f(e) \geq f(e) - b \geq f(e) - f(e) = 0.$$

Hence the capacity constraint holds whether e is a forward or backward edge.

Intuitively the algorithm will return a flow since f itself is a flow. A more formal prove involves looking at the edges depending on whether they are backward or forward edges and compare the entering flow to the exiting one.



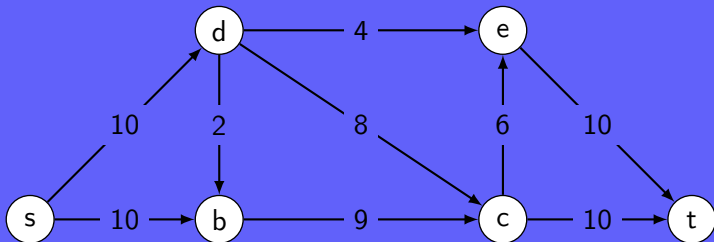
Algorithm. (*Ford-Fulkerson*)

Input : A graph G

Output: A maximum flow in G

```
1 foreach edge  $e$  in  $G$  do
2   |  $f(e) \leftarrow 0$ ;
3 end foreach
4 while there is a path  $P$  in the residual graph  $G_f$  do
5   |  $f \leftarrow \text{Augment}(f, P)$ ;
6   | update the residual graph  $G_f$ ;
7 end while
8 return  $f$ 
```

Apply Ford-Fulkerson algorithm (4.182) to our initial toy example.



In order to prove that Ford-Fulkerson algorithm (4.182) is correct we must formalize the idea of *cut* introduced in remark 4.174.

Definition

Let $G = \langle V, E \rangle$ be a flow network and s and t be the source and sink vertices, respectively.

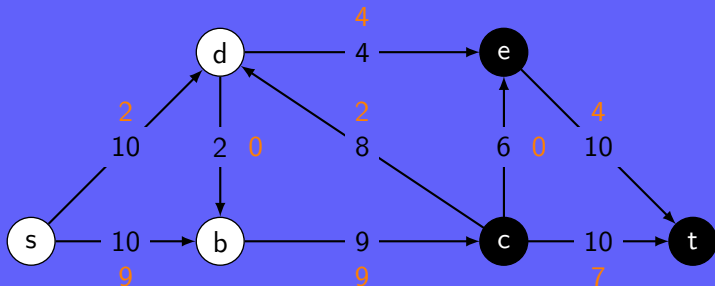
- 1 A *cut* is a partition of V into two connected subsets S and $T = V \setminus S$, with $s \in S$ and $t \in T$.
- 2 Given a flow f the *net flow* across the cut (S, T) is defined as

$$f(S, T) = \sum_{u \in S} \sum_{v \in T} f((u, v)) - \sum_{u \in S} \sum_{v \in T} f((v, u)).$$

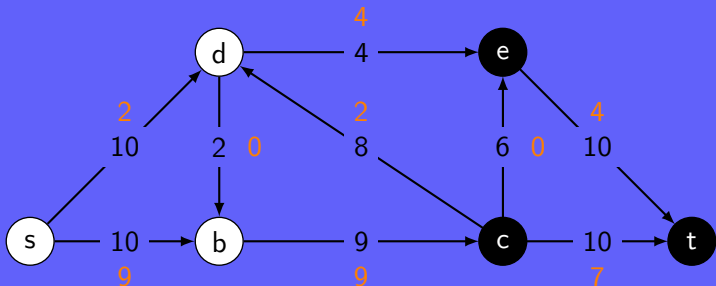
- 3 The *capacity* of the cut (S, T) is defined as

$$c(S, T) = \sum_{u \in S} \sum_{v \in T} c((u, v)).$$

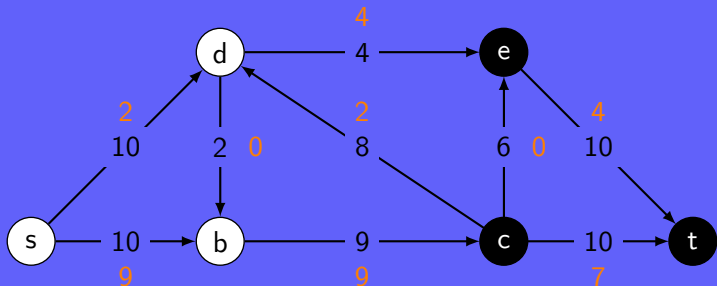
- 4 A *minimum cut* is a cut of minimum capacity over all the cuts of the network.



- Cut: $S = \{s, d, b\}$ and $T = \{c, e, t\}$
- Net flow:



- Cut: $S = \{s, d, b\}$ and $T = \{c, e, t\}$
- Net flow: $4 + 9 - 2 = 11$
- Capacity of the cut:



- Cut: $S = \{s, d, b\}$ and $T = \{c, e, t\}$
- Net flow: $4 + 9 - 2 = 11$
- Capacity of the cut: $9 + 4 = 13$

Definition

Let $G = \langle V, E \rangle$ be a flow network, s be the source, and f be a flow on G . We denote the *value* of the flow f by

$$|f| = \sum_{v \in V} f(s, v) - \sum_{v \in V} f(v, s).$$

We will now prove that given a flow f the net flow across any cut is the same.

Lemma

Let f be a flow in a flow network $G = \langle V, E \rangle$ with source s and sink t . For any cut (S, T) of G , the net flow across (S, T) is $f(S, T) = |f|$.

Proof.

From the flow conservation (4.173) for any node $u \in V \setminus \{s, t\}$

$$\sum_{v \in V} f(u, v) - \sum_{v \in V} f(v, u) = 0.$$

Adding it to the value of the flow we get

$$|f| = \sum_{v \in V} f(s, v) - \sum_{v \in V} f(v, s) + \sum_{u \in S \setminus \{s\}} \left(\sum_{v \in V} f(u, v) - \sum_{v \in V} f(v, u) \right).$$

Rewriting the right-hand sum yields

$$|f| = \sum_{v \in V} \sum_{u \in S} f(u, v) - \sum_{v \in V} \sum_{u \in S} f(v, u).$$

Proof (continued).

Noting that $S \cup T = V$ and $S \cap T = \emptyset$, the sum over the elements of V can be split over S and T to obtain

$$\begin{aligned}|f| &= \sum_{v \in S} \sum_{u \in S} f(u, v) + \sum_{v \in T} \sum_{u \in S} f(u, v) - \sum_{v \in S} \sum_{u \in S} f(v, u) - \sum_{v \in T} \sum_{u \in S} f(v, u) \\&= \sum_{v \in T} \sum_{u \in S} f(u, v) - \sum_{v \in T} \sum_{u \in S} f(v, u) + \left(\sum_{v \in S} \sum_{u \in S} f(u, v) - \sum_{v \in S} \sum_{u \in S} f(v, u) \right) \\&= \sum_{v \in T} \sum_{u \in S} f(u, v) - \sum_{v \in T} \sum_{u \in S} f(v, u) \\&= f(S, T)\end{aligned}$$

Hence the net flow across the cut (S, T) is exactly the value of the flow f . \square

Lemma

Given a flow network G and a flow f , f is upper bounded by the capacity of any cut in G .

Proof.

Using lemma 4.186 and the capacity constraint (4.173) we have

$$\begin{aligned} |f| &= f(S, T) = \sum_{v \in T} \sum_{u \in S} f(u, v) - \sum_{v \in T} \sum_{u \in S} f(v, u) \\ &\leq \sum_{v \in T} \sum_{u \in S} f(u, v) \leq \sum_{v \in T} \sum_{u \in S} c(u, v) \\ &= c(S, T) \end{aligned}$$



From the previous lemma (4.189) we can conclude that the maximum flow is upper bounded by the capacity of a minimum cut. We now prove that they are in fact equal.

Theorem (Max-flow Min-cut Theorem)

Let f be a flow in a flow network $G = \langle V, E \rangle$, with source s and sink t . Then following conditions are equivalent.

- i The flow f is a maximum flow in G .
- ii The residual network G_f contains no augmenting path.
- iii The value of the flow $|f|$ is equal to the capacity of some cut (S, T) of G .

Proof.

(i) \Rightarrow (ii) : we suppose that f is a maximum flow but G_f has an augmenting path. This generates a flow with value strictly larger than $|f|$, which contradicts lemma 4.189. ⚡

(ii) \Rightarrow (iii) : we suppose that G_f has no augmenting path, i.e. has no path from s to t . Let $S = \{v \in V \mid \text{there is a path from } s \text{ to } v\}$, and $T = V \setminus S$.

Clearly (S, T) is a cut, since $s \in S$ and $t \notin S$ for otherwise there would be a path from s to t .

Let $u \in S$ and $v \in T$ be a pair of vertices.

If $e = (u, v)$ is in E then $f(e) = c(e)$, otherwise e would be in E_f and v would be in S . ⚡

If $e = (v, u)$ is in E then $f(e) = 0$, otherwise $c_f((u, v)) = f(e)$ would be positive, (u, v) would be in E_f , and v would be in S . ⚡

Proof (continued).

For the last case, note that if neither (u, v) nor (v, u) belongs to E it means that $f((u, v)) = f((v, u)) = 0$.

Thus we have

$$\begin{aligned} f(S, T) &= \sum_{v \in T} \sum_{u \in S} f(u, v) - \sum_{v \in T} \sum_{u \in S} f(v, u) \\ &= \sum_{v \in T} \sum_{u \in S} c(u, v) - \sum_{v \in T} \sum_{u \in S} 0 \\ &= c(S, T) \end{aligned}$$

Therefore, by lemma 4.186, $|f| = c(S, T)$.

(iii) \Rightarrow (i) : combining lemma 4.189 with $|f| = c(S, T)$ yields (i).



While the Max-flow Min-cut theorem (4.190) provides a proof of accuracy for Ford-Fulkerson algorithm (4.182), the question of its efficiency remains unanswered.

In fact the algorithm still features an unclear request: “while there is a path P in the residual graph G_f ”. So the question boils down to knowing “how to determine all the paths and what is the cost”.

A simple implementation of the algorithm would complete as the value of the flow keeps increasing by at least one unit while never exceeding the maximum flow f^* . Therefore in the worst case the complexity is $\mathcal{O}(|E||f^*|)$.

On the other hand if edges have a large capacity this strategy is not optimal and finding a “good path” becomes of a major importance. In practice this can be done using the breadth-first search algorithm which will allow us to determine the shortest path.

Algorithm. (*Breadth-First Search (BFS)*)

Input : a graph G and a starting vertex s

Output: none – only access all the vertices accessible from s

```
1 foreach vertex  $v$  in  $G$  do  $v.dist \leftarrow \infty$ ;  $v.parent \leftarrow \text{NULL}$ ;  
2  $Q \leftarrow \emptyset$ ;  $s.dist \leftarrow 0$ ; append  $s$  to  $Q$ ;  
3 while  $Q \neq \emptyset$  do  
4    $v \leftarrow$  first element in  $Q$ ; remove  $v$  from  $Q$ ;  
5   foreach vertex  $u$  adjacent to  $v$  do  
6     if  $u.dist = \infty$  then  
7        $u.dist \leftarrow v.dist + 1$ ;  $u.parent \leftarrow v$ ; append  $u$  to  $Q$ ;  
8     end if  
9   end foreach  
10 end while
```

Theorem

Given a graph $G = \langle V, E \rangle$, the complexity of Breadth-First search is $\mathcal{O}(|E|)$.

Proof.

The result is straight forward as the whole adjacency list is to be scanned.



We now prove that BFS correctly computes the shortest path distance between s and any vertex v .

Lemma

Let $G = \langle V, E \rangle$ be a graph and s be a source vertex in G . Upon running BFS on G , for any vertex v , the computed value $v.dist$ is larger than the shortest-path distance between s and v .

Proof.

Let δ_v be the shortest distance between s and v . We will prove the result by induction on the number of “append” operations, the hypothesis being $v.dist \geq \delta_v$.

Base case: the first element to be appended to the list is s for which $\delta_s = s.dist = 0$; for all other vertices $\delta_v \leq v.dist = \infty$.

Proof (continued).

Induction step: let v and u be two vertices where u is discovered from v . Noting that $u.dist = v.dist + 1$ we apply the induction hypothesis to v and get

$$\begin{aligned} u.dist &\geq \delta_v + 1 \\ &\geq \delta_u \end{aligned}$$

Noting that a vertex is never appended more than once, $u.dist$ will not be updated. Therefore the induction principle applies and for any vertex v in V , $v.dist$ is larger than the shortest-path distance between s and v . □

We now want to prove that (i) nodes in Q are ordered with respect to their distance, and (ii) the first and last have a distance difference of at most one.

Lemma

Let $G = \langle V, E \rangle$ be a graph and s be a source vertex in G . Upon running BFS on G , Q contains the vertices $\{v_1, v_2, \dots, v_r\}$ where v_1 and v_r are the head and tail, respectively. Then $v_r.dist \leq v_1.dist + 1$ and for $1 \leq i \leq r - 1$, $v_i.dist \leq v_{i+1}.dist$.

Proof.

We prove the result by induction on the number of operations on Q , that is we take into account both “append” and “remove”.

Base case: clearly the result holds when Q only contains s .

Induction step: two cases must be considered, when (i) removing and (ii) appending an element.

(i) When removing v_1 two cases can arise: either Q becomes empty and in this case the result holds vacuously, or the second element, v_2 becomes the head.

Proof (continued).

If v_2 is the new head then by the induction hypothesis $v_1.dist$ was less than $v_2.dist$. But as by induction $v_r.dist$ was less than $v_1.dist + 1$, we conclude that $v_r.dist \leq v_2.dist + 1$. All the remaining equalities being unaffected by the change of head, the result holds when removing an element.

(ii) When exploring the vertices adjacent to some vertex v , v has already been removed from Q . Let u be a neighboring vertex from v . Then u is appended to Q and can be renamed as v_{r+1} .

By the induction hypothesis $v.dist \leq v_1.dist$. Thus $v_{r+1}.dist$ is the same as $v.dist + 1$ which is less than $v_1.dist + 1$. Furthermore by the induction hypothesis $v_r.dist \leq v.dist + 1$, such that $v_r.dist$ is less than $v.dist + 1 = v_{r+1}.dist$. All the remaining equalities are unaffected.

Therefore by the induction principle lemma 4.198 holds. □

Theorem

Let $G = \langle V, E \rangle$ be a graph and s be a source vertex in G . Then BFS discovers all the vertices v in V reachable from s and upon termination for all of them $v.dist$ is the shortest distance δ_v .

Proof.

Let $v \in V$ be a vertex such that $v.dist \neq \delta_v$. Then by lemma 4.196 $v.dist > \delta_v$. Clearly v must be reachable from s otherwise δ_v would be $\infty \geq v.dist$.

Let u be the vertex immediately preceding v on the shortest path to v . Then $\delta_v = \delta_u + 1$ and $u.dist = \delta_u$. Hence

$$v.dist > \delta_v = \delta_u + 1 = u.dist + 1 \quad (4.1)$$

The vertex v can be in three states: (i) neither in Q nor visited, (ii) in Q , and (iii) not in Q but visited.

Proof (continued).

(i) When v is visited $v.dist$ is set to $u.dist+1$, which contradicts equation (4.1).⚡

(ii) If v is in Q then it means it was added during the exploration of the neighbors of a vertex w . Either $u = w$ or w was removed from Q earlier than u . So $v.dist = w.dist + 1$ and $w.dist \leq u.dist$. Therefore $v.dist = w.dist + 1 \leq u.dist + 1$, contradicting (4.1).⚡

(iii) If v has already been removed from Q then clearly $v.dist < u.dist$, which once again contradicts (4.1).⚡

Hence for all v in V , $v.dist = \delta_v$. Moreover all the vertices from V must be discovered otherwise if such a vertex v existed, then δ_v would be smaller than $v.dist = \infty$. □

The Edmonds-Karp algorithm

In order to evaluate the complexity of the Ford-Fulkerson algorithm (4.182), the question of how to find an augmenting path had to be answered.

We now prove that a shortest path in the residual network can be discovered using BFS and then used as an augmenting path. The resulting algorithm is called the Edmonds-Karp algorithm.

Lemma

Let $G = \langle V, E \rangle$ be a flow network with source s and sink t . If the Edmonds-Karp algorithm is run on G then for any vertex v in $V \setminus \{s, t\}$, the shortest-path distance $\delta_{f,v}$ in the residual network G_f increases monotonically with each flow augmentation.

Shortest path as augmenting path

Proof.

Suppose that there exists a vertex $v \in V \setminus \{s, t\}$ such that a flow augmentation causes the shortest-path distance between s and v to decrease.

Let f and f' be the flows just before and after the augmentation that decreases the shortest-path distance, respectively. Then for v we have $\delta_{f',v} < \delta_{f,v}$. If u is the vertex visited just before v in $G_{f'}$ then $\delta_{f',u} = \delta_{f',v} - 1$.

From the choice of v , the shortest distance between s and u did not decrease on the flow augmentation. Therefore $\delta_{f,u} \leq \delta_{f',u}$ and (u, v) cannot belong to E_f . In fact if this was the case then we would have

$$\begin{aligned} \delta_{f,v} &\leq \delta_{f,u} + 1 \\ &\leq \delta_{f',u} + 1 \\ &= \delta_{f',v}. \end{aligned}$$

This contradicts the assumption $\delta_{f',v} < \delta_{f,v}$.

Shortest path as augmenting path

Proof (continued).

Since (u, v) does not belong to E_f but then belongs to $E_{f'}$ it means that the flow has been increased from v to u . However as the Edmonds-Karps applies BFS to augment the flow along the shortest path we conclude that the shortest path from s to u has (v, u) as its last edge. Therefore,

$$\begin{aligned}\delta_{f,v} &= \delta_{f,u} - 1 \\ &\leq \delta_{f',u} - 1 \\ &= \delta_{f',v} - 2.\end{aligned}$$

Again this contradicts the assumption $\delta_{f',v} < \delta_{f,v}$.

Hence there exists no vertex $v \in V \setminus \{s, t\}$ such that a flow augmentation causes the shortest-path distance between s and v to decrease. \square

Theorem

Let $G = \langle V, E \rangle$ be a flow network with source s and sink t . If the Edmonds-Karp algorithm is run on G then it returns a maximum flow in time $\mathcal{O}(|V||E|^2)$.

Proof.

Given a path P in the residual network G_f , we call an edge e such that $c_f(P) = c_f(e)$ a *critical edge*. We immediately notice that (i) such an edge disappears from the residual network as soon as an augmentation is performed along P , and (ii) at least one edge on any augmenting path is critical.

In order to determine the complexity of the Edmonds-Karp algorithm we must figure out how many flow augmentation are performed.

Proof (continued).

Let u and v be two vertices and $e = (u, v)$ an edge in E . Since augmenting paths are shortest paths we have $\delta_{f,v} = \delta_{f,u} + 1$.

If e is a critical edge then it will disappear from the residual network as soon as the flow is augmented. It can however reappear later after the flow from u to v is decreased, that is if (v, u) is part of an augmenting path. The question is then to know how many times it can disappear and reappear.

Let f' be the flow in G when (v, u) appears in an augmenting path. Then $\delta_{f',u} = \delta_{f',v} + 1$. Moreover as $\delta_{f,v} \leq \delta_{f',v}$ (4.202) we get

$$\begin{aligned}\delta_{f',u} &= \delta_{f',v} + 1 \\ &\geq \delta_{f,v} + 1 \\ &= \delta_{f,u} + 2\end{aligned}$$

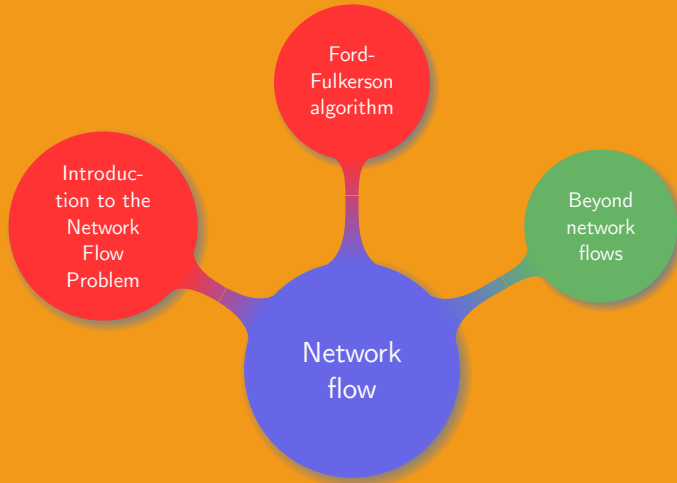
Proof (continued).

Thus, if e is a critical edge then the next time it becomes critical the distance between the s and u has increased by at least 2, while it originally was at least 1.

Then observing that the path from s to u cannot contain u , or t , we conclude that the distance will be at most $|V| - 2$.

Finally by combining the two bounds we conclude that an edge is critical no more than $|V|/2$ times. And as the number of edges in the residual network is $\mathcal{O}(|E|)$, the total number of critical edges during the execution of the algorithm is $\mathcal{O}(|V||E|)$.

Recalling that an augmenting path has at least one critical edge, we loop $\mathcal{O}(|V||E|)$ times in the Edmonds-Karp algorithm. For each of them, BFS with complexity $\mathcal{O}(|E|)$ (theorem 4.195) is run, leading to a total cost of $\mathcal{O}(|V||E|^2)$. \square



When studying computability theory the importance of finding similarities between problems was highlighted (2.131). The idea behind this approach is to solve new problems by deriving appropriate algorithms from known ones.

The difficulty is therefore to view a problem from a different perspective. In practice this is similar to determining polynomial reductions in computability theory (2.114): one wants to efficiently “rephrase ” a given problem into an new one for which a solution is known.

We now study such an example as we solve the maximum bipartite matching problem using network flows.

Definition

Let $G = \langle V, E \rangle$ be a graph.

- ① If $V = L \cup R$, with L and R two disjoint sets, every vertex in V has at least one incident edge, and for any edge (u, v) either $u \in L$ and $v \in R$ or $u \in R$ and $v \in L$, then G is called *bipartite graph*.
- ② A *matching* is a subset M of E such that for any vertex $v \in V$ at most one edge in M is incident to v .
- ③ A *maximum matching* is a matching of maximum cardinality.

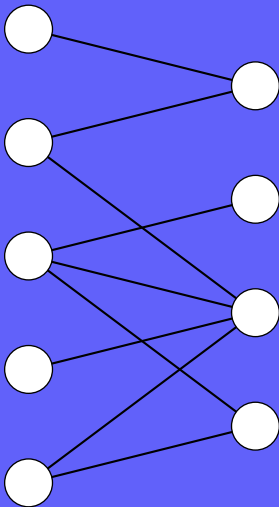
Definition

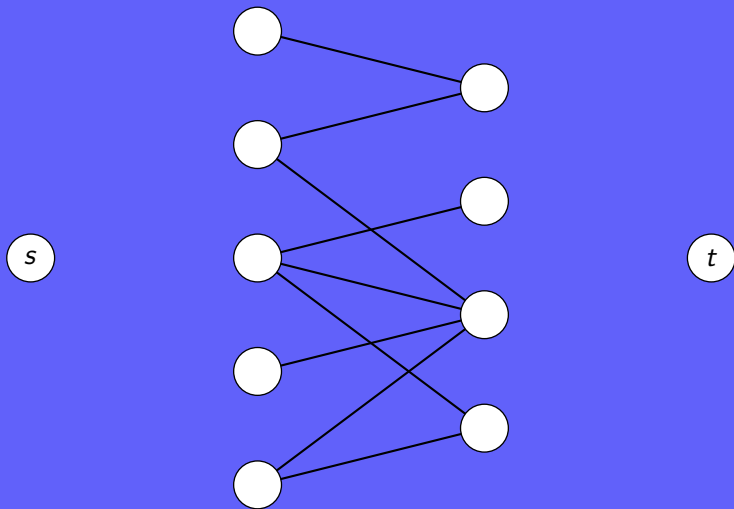
Let $G = \langle V, E \rangle$ be a graph.

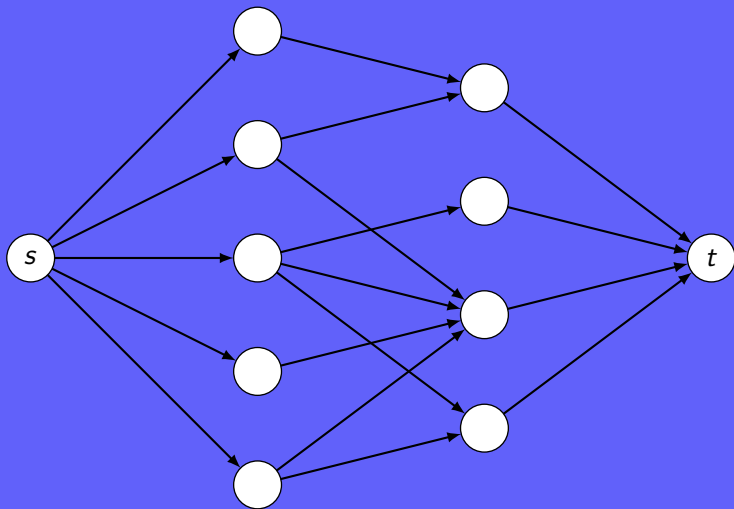
- 1 If $V = L \cup R$, with L and R two disjoint sets, every vertex in V has at least one incident edge, and for any edge (u, v) either $u \in L$ and $v \in R$ or $u \in R$ and $v \in L$, then G is called *bipartite graph*.
- 2 A *matching* is a subset M of E such that for any vertex $v \in V$ at most one edge in M is incident to v .
- 3 A *maximum matching* is a matching of maximum cardinality.

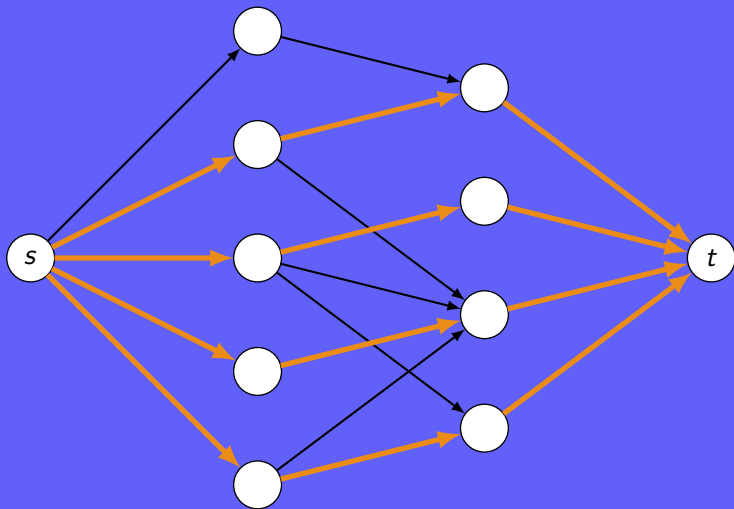
Problem (Maximum Bipartite Matching Problem)

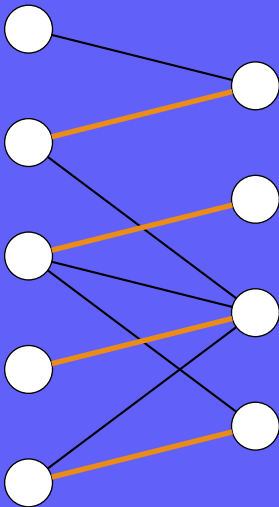
Given a bipartite graph, determine a maximum matching.











To be able to apply Ford-Fulkerson method we formalize the transformation of a bipartite graph $G = \langle V, E \rangle$ into a network flow.

First define $G' = \langle V', E' \rangle$, with $V' = V \cup \{s, t\}$, s and t being the source and the sink, respectively. Since G is a bipartite graph V can be partitioned into two sets L and R , and

$$E' = \{(s, u) / u \in L\} \cup \{(u, v) / (u, v) \in E\} \cup \{(v, t) / v \in R\}.$$

It then suffices to assign a unit capacity to each edge in E' .

To be able to apply Ford-Fulkerson method we formalize the transformation of a bipartite graph $G = \langle V, E \rangle$ into a network flow.

First define $G' = \langle V', E' \rangle$, with $V' = V \cup \{s, t\}$, s and t being the source and the sink, respectively. Since G is a bipartite graph V can be partitioned into two sets L and R , and

$$E' = \{(s, u) / u \in L\} \cup \{(u, v) / (u, v) \in E\} \cup \{(v, t) / v \in R\}.$$

It then suffices to assign a unit capacity to each edge in E' .

At this stage an important step, that should not be omitted, is to check the cost of the transformation.

Since each vertex in V has at least one incident edge, $|E|$ is greater or equal to $|V|/2$. Thus the number of edges in G is smaller than the one in G' and $|E'| = |E| + |V| \leq 3|E|$. Hence $|E'| = \Theta(|E|)$, meaning that the transformation can be efficiently performed.

Lemma

Let $G = \langle V, E \rangle$ be a bipartite graph with vertex partition $V = L \cup R$, and $G' = \langle V', E' \rangle$ be its corresponding flow network. There is a matching M in G if and only if there exists a flow f in G' . In particular the value of the flow $|f|$ is equal to the cardinality of the matching $|M|$.

Proof.

Let M be a matching in G . For an edge $e = (u, v)$ in E' , define a flow f as $f(s, u) = f(u, v) = f(v, t) = 1$ if e is in M and 0 otherwise. Clearly f satisfies both the capacity constraint and the flow conservation properties (4.173).

Moreover as G is a bipartite graph a simple cut can be defined as $(S, T) = (L \cup \{s\}, R \cup \{t\})$. Observing that the net flow across the cut (S, T) is equal to $|M|$, we apply lemma 4.186 and get $|f| = |M|$.

Proof (continued).

To prove the converse define a flow f on G' and M to be

$$M = \{(u, v) / u \in L, v \in R, \text{ and } f(u, v) > 0\}.$$

By construction each vertex u in L has a single entering edge (s, u) with capacity at most 1. Therefore by the flow conservation property no more than 1 unit can leave on at most one edge. Thus a unit can enter u if and only if there is at most one v in R such that (u, v) is in M . Hence M is a matching.

Consequently for any matched vertex u in L , $f(s, u) = 1$, while $f(u, v) = 0$ for any edge in $E \setminus M$. This means that the net flow across cut the $(L \cup \{s\}, R \cup \{t\})$ is equal to $|M|$, which is exactly $|f|$, by lemma 4.186. \square

Remark.

Over this whole chapter the net flow was implicitly assumed to only take integral values. However in the previous proof (4.213) it was explicitly mentioned that a unit flow cannot be split, which is true only in the case of integer valued flows.

In particular if the capacity function only takes integral values then the Ford-Fulkerson method returns an integer valued flow.

Remark.

Over this whole chapter the net flow was implicitly assumed to only take integral values. However in the previous proof (4.213) it was explicitly mentioned that a unit flow cannot be split, which is true only in the case of integer valued flows.

In particular if the capacity function only takes integral values then the Ford-Fulkerson method returns an integer valued flow.

Theorem

Let $G = \langle V, E \rangle$ be a bipartite graph. The cardinality of a maximum matching M is equal to the value of the maximum flow in the flow network G' corresponding to G . This maximum matching is determined in time $\mathcal{O}(|V||E|)$.

Proof.

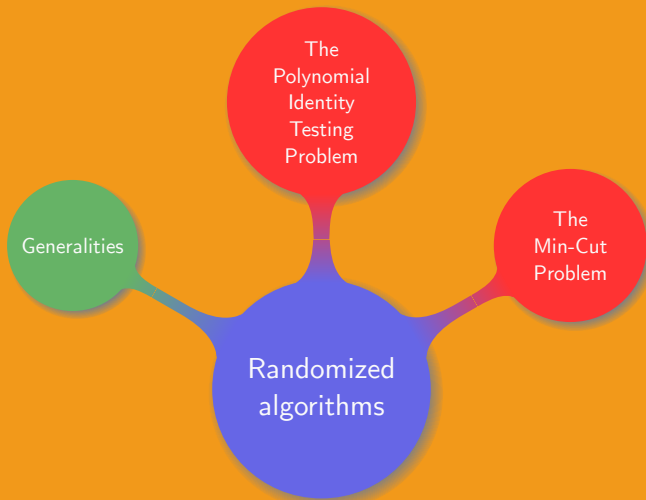
Suppose M is a maximum matching in G while the corresponding flow in G' is not maximum. Then we can find a flow f' such that $|f'| > |f|$. But as noted in remark 4.215, both f and f' take integral values and there exists a matching M' corresponding to f' . Therefore we get

$$|M| = |f| < |f'| = |M'| \quad \text{.}$$

Similarly if f is a maximum flow in G' , then M is a maximum matching in G .

As any matching in a bipartite graph has cardinality at most $\min(|L|, |R|) = \mathcal{O}(|V|)$ the value of the maximum flow in G' is $\mathcal{O}(|V|)$. Therefore by 4.193 and 4.212 a maximum matching can be found in time $\mathcal{O}(|V||E'|) = \mathcal{O}(|V||E|)$. \square

5. Randomized algorithms



Along chapter 2, section 3 (2.103) we studied:

- How to model computers
- What an algorithm is
- Various classes of problems

Along chapter 2, section 3 (2.103) we studied:

- How to model computers
- What an algorithm is
- Various classes of problems

Computation on a deterministic Turing machine depends on:

- The state of the machine
- The symbol read on the tape

For any situation at most one action is performed

Computation on a non-deterministic Turing machine:

- The machine has a state and a symbol is read on the tape
- The machine branches into many copies
- The machine transitions into one of the copies

When the machine is run more than once:

- Different paths are chosen
- Computation cannot be exactly reproduced

Non-deterministic Turing machine

Computation on a non-deterministic Turing machine:

- The machine has a state and a symbol is read on the tape
- The machine branches into many copies
- The machine transitions into one of the copies

When the machine is run more than once:

- Different paths are chosen
- Computation cannot be exactly reproduced

Remark.

A non-deterministic Turing machine can be simulated by a deterministic Turing machine with three tapes.

Two main types of algorithms:

- Deterministic: follow a fixed sequence of steps; output is completely determined by the input
- Probabilistic: add some randomness to the process
 - Monte Carlo algorithm: returns either True or unknown; increasing the number of test decreases the probability of having a “false positive”
 - Las Vegas algorithms: always returns a correct result; the running time is however random; the time complexity cannot be precisely evaluated

Quick sort:

- Worst case: $\mathcal{O}(n^2)$ vs. average case $\mathcal{O}(n \log n)$
- Fixed pivot: if the list to sort is originally structured worst complexity is likely to apply
- Random pivot: even if the list to sort is originally structured choosing a random pivot is equivalent to sorting with a fixed pivot on a randomly ordered list
- Running time: random depending on the choice of the pivot

Primality testing:

- AKS: deterministic in time $\tilde{O}(n \log^{12} n)$
- Miller-Rabbin:
 - Complexity: $\tilde{O}(k \log^2 n)$, with k the number of witnesses
 - Returns composite or unknown
 - Unknown means prime with probability $1 - 4^{-k}$

Definitions

- ① A non-deterministic Turing machine which chooses a random transitions according to some probability distribution is called a *probabilistic Turing machine*.
- ② A language L is in the Bounded-error Probabilistic Polynomial time complexity class (BPP) if and only if there is a probabilistic Turing machine M such that
 - For any input, M runs in polynomial time;
 - For all x in L , M returns 1 with probability larger than $1/2 + \epsilon$, $\epsilon > 0$;
 - For all x not in L , M returns 1 with probability less than $1/2 - \epsilon$, $\epsilon > 0$;

Remark.

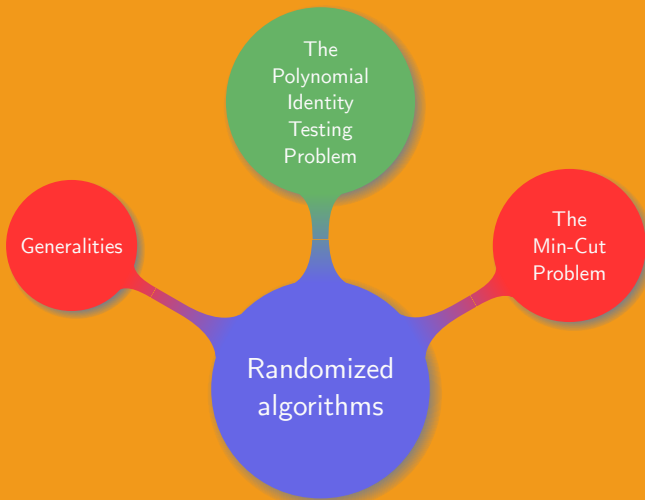
In practice instead of taking $1/2 + \epsilon$, $2/3$ and $1/3$ are often chosen.

Reasons for using probabilistic algorithms:

- No polynomial time deterministic algorithm is known
- Available polynomial time deterministic algorithms are slower

Open questions:

- Is $\mathcal{P} = \text{BPP}$?
- Are probabilistic algorithm more powerful than deterministic ones?
- What is the relative power of probabilistic and deterministic computations?



The polynomial identity testing problem

Problem (Polynomial Identity Testing (PIT))

Let P be a multivariate polynomial over some field. Decide whether P is identically zero.

Remark.

Identically zero means that after expanding P it reduces to zero. In particular this is a different problem from the Evaluate to Zero Everywhere (EZE) problem where one wants to decide if an n -multivariate polynomial evaluates to zero for all x_i , $1 \leq i \leq n$.

The univariate case can be easily solved:

- The polynomial is given as a sum of monomials:
 - Reduce the sum
 - Check all the monomials and return True if they all equal zero
- The polynomial P of degree d is in a more complex form:
 - Arbitrarily choose $d + 1$ points (e.g. $0, \dots, d$)
 - Evaluate P at those $d + 1$ points
 - Conclude that $P = 0$ if and only if they all evaluate to zero

The univariate case can be easily solved:

- The polynomial is given as a sum of monomials:
 - Reduce the sum
 - Check all the monomials and return True if they all equal zero
- The polynomial P of degree d is in a more complex form:
 - Arbitrarily choose $d + 1$ points (e.g. $0, \dots, d$)
 - Evaluate P at those $d + 1$ points
 - Conclude that $P = 0$ if and only if they all evaluate to zero

Why is this not working anymore in the case of a multivariate polynomial?

Definition

- ① A *monomial* is an expression of the form $\alpha \prod_{i=1}^n x_i^{\beta_i}$, where α is an element from a base field, the x_i are n variables, and the β_i are positive integers.
- ② The *total degree of a monomial* is $\sum_i \beta_i$.
- ③ The *total degree of a polynomial* is the largest total degree among all the monomials composing the polynomial.

Definition

- 1 A *monomial* is an expression of the form $\alpha \prod_{i=1}^n x_i^{\beta_i}$, where α is an element from a base field, the x_i are n variables, and the β_i are positive integers.
- 2 The *total degree of a monomial* is $\sum_i \beta_i$.
- 3 The *total degree of a polynomial* is the largest total degree among all the monomials composing the polynomial.

Lemma (Schwartz-Zippel)

Let P be a n -multivariate polynomial of total degree d , that is not identically zero, over a field \mathbb{F} . For y_1, \dots, y_n , chosen uniformly and independently from a finite set $S \subset \mathbb{F}$,

$$\Pr [P(y_1, \dots, y_n) = 0] \leq \frac{d}{|S|}.$$

Proof.

We proceed by induction on the number of variables n .

Base case: For $n = 1$ the result is clear as a univariate polynomial of degree d has at most d roots.

Induction step: we assume that the result is true for an $(n - 1)$ -multivariate polynomial and prove it is also true in the case of n variables. Let k be the largest power of X_1 in any monomial composing P . Then

$$P(X_1, \dots, X_n) = \sum_{i=0}^k X_1^i Q_i(X_2, \dots, X_n).$$

By construction, Q_k is not identically zero, its total degree is at most $d - k$, and it has $n - 1$ variables. Therefore by the induction hypothesis we get

$$\Pr [Q_k(y_2, \dots, y_n) = 0] \leq \frac{d - k}{|S|}.$$

Proof (continued).

For y_2, \dots, y_n in \mathbb{F} , we call \mathcal{E}_1 the event $Q_k(y_2, \dots, y_n) = 0$.

Selecting y_2, \dots, y_n such that \mathcal{E}_1 does not occur, we define $R(X_1)$ to be the polynomial

$$R(X_1) = \sum_{i=0}^k X_1^i Q_i(y_2, \dots, y_n) = P(X_1, y_2, \dots, y_n).$$

Clearly $R(X_1)$ is not identically zero since \mathcal{E}_1 did not occur, meaning that X_1^k has a non zero coefficient. Therefore

$$\Pr [R(y_1) = 0 \mid \neg \mathcal{E}_1] \leq \frac{k}{|S|}.$$

Let \mathcal{E}_2 be the event $R(y_1) = 0$, which can also be stated as $P(y_1, \dots, y_n) = 0$. In order to prove the lemma it remains to bound $\Pr [\mathcal{E}_2]$.

Proof (continued).

As we have already bounded $\Pr[\mathcal{E}_2 \mid \neg \mathcal{E}_1]$ and $\Pr[\mathcal{E}_1]$ we can rewrite $\Pr[\mathcal{E}_2]$ as

$$\begin{aligned}\Pr[\mathcal{E}_2] &= \Pr[\mathcal{E}_2 \wedge \mathcal{E}_1] + \Pr[\mathcal{E}_2 \wedge \neg \mathcal{E}_1] \\ &= \Pr[\mathcal{E}_2 \wedge \mathcal{E}_1] + \Pr[\mathcal{E}_2 \mid \neg \mathcal{E}_1] \Pr[\neg \mathcal{E}_1] \\ &\leq \Pr[\mathcal{E}_1] + \Pr[\mathcal{E}_2 \mid \neg \mathcal{E}_1] \\ &\leq \frac{d-k}{|S|} + \frac{k}{|S|} = \frac{d}{|S|}\end{aligned}$$

Hence, by the induction principle Schwartz-Zippel lemma holds.



Remark.

Schwartz-Zippel lemma (5.230) says that when evaluating P at a random point there is a very low probability of finding a root. It however does not mean that a polynomial over \mathbb{R} has finitely many roots. For instance $P(X_1, X_2) = X_1$ has infinitely many roots.

Algorithm.

Input : $P(X_1, \dots, X_n)$ of degree d and a field \mathbb{F} , l and k

Output: not zero or probably zero

- 1 Select a subset $S \subset \mathbb{F}$ of size $\geq ld$;
 - 2 **for** $i \leftarrow 1$ **to** k **do**
 - 3 $(y_1, \dots, y_n) \leftarrow \text{rand}(S)$;
 - 4 **if** $P(y_1, \dots, y_n) \neq 0$ **then return** not zero;
 - 5 **end for**
 - 6 **return** probably zero;
-

Algorithm.

Input : $P(X_1, \dots, X_n)$ of degree d and a field \mathbb{F} , l and k

Output: not zero or probably zero

- 1 Select a subset $S \subset \mathbb{F}$ of size $\geq ld$;
 - 2 **for** $i \leftarrow 1$ **to** k **do**
 - 3 $(y_1, \dots, y_n) \leftarrow \text{rand}(S)$;
 - 4 **if** $P(y_1, \dots, y_n) \neq 0$ **then return** not zero;
 - 5 **end for**
 - 6 **return** probably zero;
-

Theorem

Let $P(X_1, \dots, X_n)$ be a polynomial. If $P(X_1, \dots, X_n)$ is not identically zero, then algorithm 5.234 returns the correct result with probability at least $1 - 1/l^k$.

Proof.

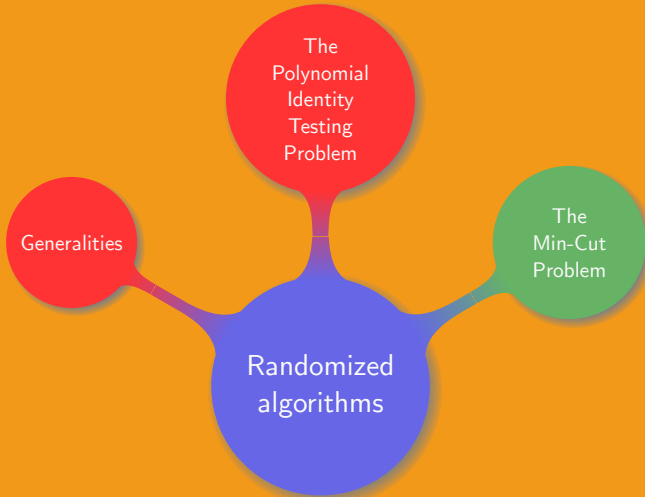
For each iteration of the loop the probability of testing a root is at most $1/l$. Each choice of element in S is independent of the previous ones, such that the probability of passing all the tests while $P(X_1, \dots, X_n)$ is not zero is at most $1/l^k$.

Hence the probability of returning “probably zero” while $P(X_1, \dots, X_n)$ is not identically zero is at least $1 - 1/l^k$.



Remark.

If the field \mathbb{F} has less elements than the number of roots of $P(X_1, \dots, X_n)$ then Schwartz-Zippel lemma (5.230) is of no use. It is however possible to overcome this problem by applying algorithm 5.234 to the polynomial $P(X_1, \dots, X_n)$ over an extension field of \mathbb{F} .



Problem (Min-Cut)

Given a connected multi-graph G , determine the minimum number of edges that must be removed such that G becomes disconnected.

Common applications:

- Split a problem for parallel programming
- Optimize a divide and conquer strategy
- Segment an image into regions of similar color/texture

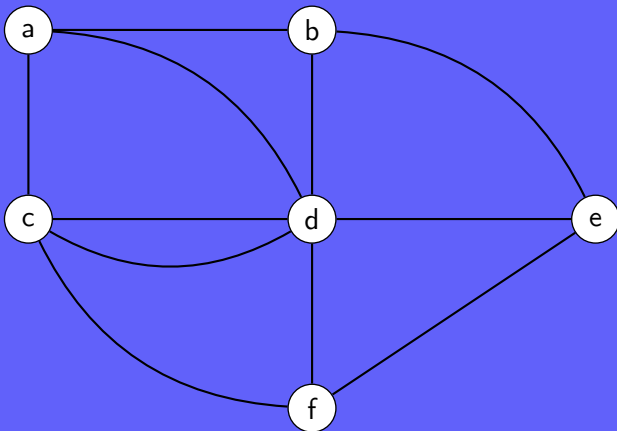
Algorithm. (*Karger*)

Input : a graph $G = \langle V, E \rangle$

Output : an upper bound on the min-cut

```
1 Function Contract( $G, e$ ):  
2   |   remove all edges between  $e.s$  and  $e.t$ ;  
3   |   foreach edge  $e'$  containing  $e.s$  do  
4   |   |   disconnect  $e'$  from  $e.s$  and reconnect it to  $e.t$ ;  
5   |   end foreach  
6   |    $S_{e.t} \leftarrow S_{e.t} \cup S_{e.s}$ ; remove vertex  $e.s$  from  $V$ ;  
7   |   return  $G$ ;  
8 end  
9 foreach  $v \in G.V$  do  $S_v = \{v\}$ ;  
10 while  $G$  has more than two vertices do  
11 |    $e \leftarrow \text{rand}(G.E)$ ;  $G \leftarrow \text{Contract}(G, e)$ ;      /* edge  $e = (e.s, e.t)$  */  
12 end while  
13 return  $|G.E|$ ;
```

Apply Karger's algorithm to the following graph G and find two different final values for $|G.E|$.



The process described in the function `Contract` of Karger's algorithm (5.238) is called *edge contraction*.

As observed in the previous exercise (5.239) running Karger's algorithm can lead to various values. The question is then to figure out how to retrieve the right one with high probability. In order to do this we first evaluate the cost of a contraction.

Lemma

Let $G = \langle V, E \rangle$ be a multigraph. A single contraction in Karger's algorithm takes $\mathcal{O}(|V|^2)$ time.

Proof.

Assuming G is represented using a linked list at each vertex, a contraction consists in merging two lists, while ensuring no loop subsists. Therefore the list is of size $\mathcal{O}(|V|^2)$. \square

Lemma

Let $G = \langle V, E \rangle$ be a multigraph, and (S, T) be a min-cut. Then

$$\Pr [\text{Karger's algorithm ends with } (S, T)] \geq \frac{1}{\binom{|V|}{2}}.$$

Proof.

First observe that if an edge (u, v) is contracted then only the cuts containing both u and v are unchanged. Therefore for Karger's algorithm to succeed the minimum cut (S, T) must remain untouched over all the random edge selections. Denoting the number of vertices by n , $n - 2$ contractions are to be performed in order to have only two vertices left. Naming those edges $\{e_1, \dots, e_{n-2}\}$, the goal is to determine the probability to choose a proper edge at each iteration of the algorithm.

Proof (continued).

Let k denote the size of a minimum cut in G . Then the minimum degree of any vertex in G is at least k , otherwise a cut of size less than k could be exhibited by disconnecting a vertex of degree less than k . Therefore G has at least $nk/2$ edges.

After a contraction the new graph has one less vertex but the degree of all the vertices remains at least k . So after i steps there are $n - i$ vertices and at least $(n - i)k/2$ edges.

Since any vertex has degree at least k the probability of selecting a “bad edge”, assuming none has been chosen before, is $2/(n - i)$.

Hence we can determine the probability of never choosing a “bad edge” during the whole process and thus end with (S, T) .

Proof (continued).

This probability is given by

$$\begin{aligned}\Pr [\text{find } (S, T)] &= \Pr [e_1, \dots, e_{n-2} \notin (S, T)] \\&= \Pr [e_1 \notin (S, T)] \prod_{i=1}^{n-3} \Pr [e_{i+1} \notin (S, T) \mid e_1, \dots, e_i \notin (S, T)] \\&\geq \prod_{i=0}^{n-3} \left(1 - \frac{2}{n-i}\right) \\&= \frac{2}{n(n-1)} = \frac{1}{\binom{n}{2}}.\end{aligned}$$



Probability of finding a min-cut

Since the probability of finding a min-cut is least $1/\binom{n}{2}$ it suffices to run the algorithm $l\binom{n}{2}$, for some value l . The probability of a run to succeed is then at least

$$1 - \left(1 - \frac{1}{\binom{n}{2}}\right)^{l\binom{n}{2}} \geq 1 - e^{-l}.$$

Therefore an appropriate choice for l is $c \ln n$, which leads to an error probability of at most $1/n^c$. Hence the total running time of Karger's algorithm is $\mathcal{O}(n^4 \log n)$.

Although this approach works well, it remains slow. The main reason is the random choice of the edge: at the beginning the multigraph features many edges and the probability of selecting an edge from a minimum cut is low. However as the process advances the probability of contracting an edge in the minimum cut grows.

Algorithm. (*Karger-Stein*)

Input : a graph $G = \langle V, E \rangle$

Output : a mini-cut in G

```

1 Function FastCut( $G$ ):
2   if  $|G.V| > 6$  then
3      $G_1 \leftarrow G$ ;  $G_2 \leftarrow G$ ;
4      $t \leftarrow 1 + \frac{|G.V|}{\sqrt{2}}$ ;
5     while  $|G_1.V| \geq t$  do  $e \leftarrow \text{rand}(G_1.E)$ ;  $G_1 \leftarrow \text{Contract}(G_1, e)$ ;
6     while  $|G_2.V| \geq t$  do  $e \leftarrow \text{rand}(G_2.E)$ ;  $G_2 \leftarrow \text{Contract}(G_2, e)$ ;
7     FastCut( $G_1$ ); FastCut( $G_2$ );
8   else
9     find the min-cut by enumeration
10  end if
11  return  $|G_1.E| \leq |G_2.E| ? |G_1.E| : |G_2.E|$ 
12 end

```

Theorem

Given a multigraph with n vertices, Karger-Stein's algorithm discovers a minimum cut in time $\mathcal{O}(n^2 \log^3 n)$, with high probability.

Sketch of proof.

First note that $6 \ll |V|$ and as such finding a minimum cut by enumeration only impacts the final complexity by a constant factor.

Given a cut, observe that the probability that it survives down to t vertices is at least $\binom{t}{2} / \binom{n}{2}$. Thus for $t = n/\sqrt{2}$ the probability of success is larger than $1/2$.

Since Karger-Stein's algorithm follows a divide and conquer strategy, its complexity can be expressed by a recurrence relation.

Sketch of proof (continued).

Then recalling that a single edge contraction costs $\mathcal{O}(n^2)$ (lemma 5.240) we get

$$T(n) = 2 \left(n^2 + T \left(\frac{n}{\sqrt{2}} \right) \right).$$

Hence by the master theorem (2.95) we conclude that the running time of Karger-Stein's algorithm is $\mathcal{O}(n^2 \log n)$.

We now consider the success probability. First as we start with n vertices and go down to $t = n/\sqrt{2}$, the success probability is $(t/n)^2 \approx 1/2$. Then at the next recursion level the graph shrinks from $n/\sqrt{2}$ to $n/2$ vertices, which means an overall success probability of about $1/4$.

Sketch of proof (continued).

More generally assume the minimum cut to still be in the graph and let $P(t)$ be the probability that a call to the algorithm with t vertices successfully computes it. Then G_i , $1 \leq i \leq 2$ still contains it with probability larger than a half. Therefore the probability that a recursive call succeeds is $1/2P(t/\sqrt{2})$. And since two recursive call are performed

$$P(t) = 1 - \left(1 - \frac{1}{2}P\left(\frac{t}{\sqrt{2}}\right)\right)^2.$$

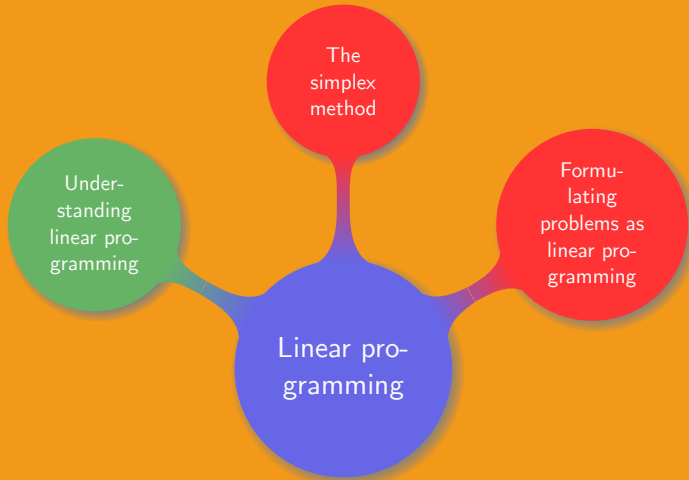
Solving this recurrence relation yields $P(n) = \Omega(1/\log n)$ ¹. This means that Karger-Stein's algorithm needs to be run about $\log^2 n$ times in order to have an error probability of at most $\mathcal{O}(1/n)$ of preserving the minimum cut. This gives a final complexity of $\mathcal{O}(n^2 \log^3 n)$. \square

¹This result is proven in the homework.

Randomized algorithms:

- Bring much flexibility compared to deterministic ones
- Are often faster than deterministic ones
- Introduce imprecision on the output or on the complexity
- Require a good knowledge of probability theory
- Have proof that are often complex, even if the algorithm can be simply expressed

6. Linear programming



A new type of problem:

- Given: limited resources or competing constraints
- Goal: maximize/minimize an objective

Problem (Linear Programming Problem)

Given a mathematical model where an *objective* is represented as a linear function and some *constraints* are expressed as equalities or inequalities, maximize the objective while respecting the constraints.

Linear programming defines a method to solve this problem.

Common applications

- Business: maximize the profit when producing various items sold at different prices
- Politics: what to say such as getting the more votes
- Airlines: following the regulation on what crew can do, maximize the number of flights possible given the number of crew members
- Exploration: decide where to drill for oil in order to maximize profit
- Diet: select food such as having enough nutriment while minimizing the cost

Three main types of algorithm exist:

- Ellipsoid algorithm: polynomial time complexity but slow in practice
- Interior-point methods: polynomial time complexity and fast in practice for large input
- Simplex algorithm: exponential time complexity but very fast in practice

Definition (Standard form)

A linear program can be expressed in *standard form*. Given n real numbers c_1, \dots, c_n , m real numbers b_1, \dots, b_m , and mn real numbers a_{ij} , $1 \leq i \leq m$, $1 \leq j \leq n$, find n real numbers x_1, \dots, x_n that maximize $\sum_{j=1}^n c_j x_j$, while being subject to

$$\begin{cases} \sum_{j=1}^n a_{ij} x_j \leq b_i, & 1 \leq i \leq m; \\ x_j \geq 0, & 1 \leq j \leq n. \end{cases} \quad (6.1)$$

In standard form all the constraints (eq. (6.1)) are defined as inequalities. Although a linear program does not necessarily requires non negative constrains, the standard form does.

Express the standard form (6.256) in a more compact way:

- Represent the a_{ij} as an $m \times n$ matrix A
- Represent the b_i as an m -vector b
- Represent the c_j as an n -vector c
- Represent the x_j as an n -vector x

In this form a linear program consists in maximizing $c^T x$, knowing that x is subject to $Ax \leq b$ and $x \geq 0$.

Linear programming and standard form

Problem

- Minimizing vs. maximizing
- Non-negative constraint on x_j
- Equality constrain $x_j = l$
- Greater than constrain $x_j \geq l$

Solution

- Negate the coefficients of the objective function
- Write $x_j = x'_j - x''_j$ and add constraints $x'_j, x''_j \geq 0$
- Write $x_j \leq l$ and $x_j \geq l$
- Multiply the inequality by -1

Linear programming and standard form

Example.

Express the following linear program into standard form. Minimize $-2x_1 + 3x_2$, subject to

$$\begin{cases} x_1 + x_2 & = 7 \\ x_1 - 2x_2 & \leq 4 \\ x_1 & \geq 0. \end{cases} \quad (6.2)$$

By negating the coefficients from the objective this is equivalent to maximizing $2x_1 - 3x_2$ while keeping the constraints (6.2).

As there is no constraint on x_2 we rewrite it $x_2 = x'_2 - x''_2$ and the problem becomes: maximize $2x_1 - 3x'_2 + 3x''_2$, subject to

$$\begin{cases} x_1 + x'_2 - x''_2 & = 7 \\ x_1 - 2x'_2 + 2x''_2 & \leq 4 \\ x_1, x'_2, x''_2 & \geq 0. \end{cases}$$

Linear programming and standard form

Finally the first constraint is rewritten as two “less or equal” inequalities such that the initial linear program can be expressed in the following standard form, with variables x_1, x_2 and x_3 in place of x_1, x'_2 and x''_2 , respectively.

Maximize $2x_1 - 3x_2 + 3x_3$, subject to

$$\left\{ \begin{array}{ll} x_1 + x_2 - x_3 & \leq 7 \\ -x_1 - x_2 + x_3 & \leq -7 \\ x_1 - 2x_2 + 2x_3 & \leq 4 \\ x_1, x_2, x_3 & \geq 0. \end{array} \right. \quad (6.3)$$

The standard form allows to define a general formalization for any linear program. However to solve such problem it is more desirable to deal with equalities rather inequalities.

Definition (Slack form)

A linear program is transformed into *slack form* by rewriting all the inequalities, but the non-negativity constraints which is omitted, as equalities and using the variable z to denote the value of the objective function.

Let $\sum_{j=1}^n a_{ij}x_j \leq b_i$ be an inequality constraint. Define a *slack variable* s and rewrite the previous constrain as

$$\begin{cases} s &= b_i - \sum_{j=1}^n a_{ij}x_j, \\ s &\geq 0. \end{cases}$$

In the new set of equalities *non-basic variables* are on the right and appear in the objective function, while *basic variables* are on the left.

For consistency reasons, when slack variables are added their name usually starts at x_{n+1} instead of s .

Example.

We carry on with example 6.259, now transforming the standard form into slack form.

From constraints (6.3) three inequalities should be rewritten as equalities. Therefore three slack variables x_4 , x_5 , and x_6 are added. Moreover the terms “maximize” and “subject to”, as well as the non-negativity constraints can be omitted. Hence the initial linear problem can be expressed in slack form as

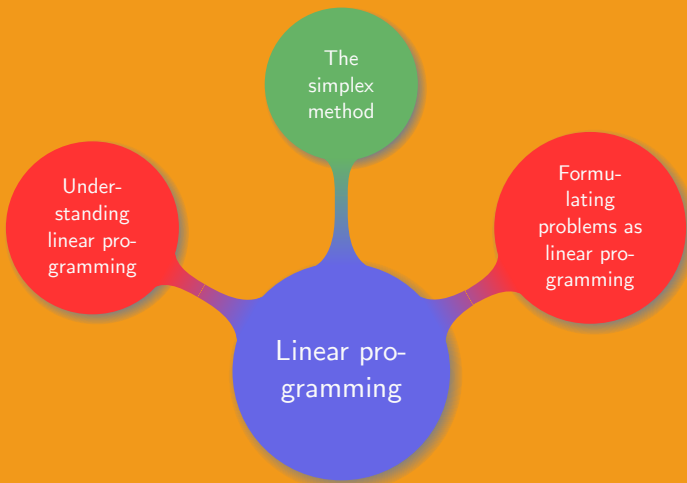
$$\begin{cases} z = & 2x_1 - 3x_2 + & x_3 \\ x_4 = & 7 - & x_1 - & x_2 + & x_3 \\ x_5 = & -7 + & x_1 + & x_2 - & x_3 \\ x_6 = & 4 - & x_1 + 2x_2 - 2x_3. \end{cases}$$

In a more compact definition a linear program can be expressed in its slack form as

$$\begin{cases} z = v + \sum_{j \in N} c_j x_j, \\ x_i = b_i - \sum_{j \in N} a_{ij} x_j, \quad \text{for } i \in B, \end{cases}$$

where z , x_j , b_i , a_{ij} , and c_j are as previously defined, while N defines the set of the non-basic variables, B represents the indices for the set of the basic variables, and v is an optional constant term in the objective function.

Therefore, in terms of algorithm a slack form can be represented by a data structure composed of the fields N, B, A, b, c, v .



Given a linear program in slack form the goal is to solve it while maximizing the value z . In fact the system of equations provided by the constraints is very similar to a regular system of equations, which can usually be solved by *Gaussian elimination*.

However note that a slack system usually features more variables than equations. Therefore there might be an infinite number of solutions. A solution such that $x_j \geq 0$ for all j is called a *feasible solution*.

A *basic solution* is found when setting all the non-basic variables to 0 on the right hand while computing the value of the basic variables on the left hand side. Assuming this basic solution is feasible the goal is to reformulate the linear program such as to maximize z .

The value z is maximized by selecting a non-basic variable x_k and increasing it as much as possible without violating any of the constraints. Then x_k is transferred on the left side to become a basic variable and the previous basic variable x_l is moved to the right. All other constraints, as well as the objective function are updated.

Example.

Given the following linear problem in slack form determine x_1 , x_2 , and x_3 such as to maximize z .

$$\begin{cases} z = & 3x_1 + x_2 + 2x_3 \\ x_4 = 30 - x_1 - x_2 - 3x_3 \\ x_5 = 24 - 2x_1 - 2x_2 - 5x_3 \\ x_6 = 36 - 4x_1 - x_2 - 2x_3. \end{cases}$$

The basic solution is $(x_1, \dots, x_6) = (0, 0, 0, 30, 24, 36)$. Then the “fastest” way to increase is to use x_1 . Looking at the constraints we observe that augmenting x_1 decreases x_4 , x_5 , and x_6 . However because of the non-negativity constraint none of them can become negative. Therefore x_1 should not be increased by more than 9 and

$$x_1 = 9 - \frac{x_2}{4} - \frac{x_3}{2} - \frac{x_6}{4}.$$

As x_1 is now a basic variable all the constraints and the objective are rewritten

$$\begin{cases} z = 27 + \frac{x_2}{4} + \frac{x_3}{2} - \frac{3x_6}{4} \\ x_1 = 9 - \frac{x_2}{4} - \frac{x_3}{2} - \frac{x_6}{4} \\ x_4 = 21 - \frac{3x_2}{4} - \frac{5x_3}{2} + \frac{x_6}{4} \\ x_5 = 6 - \frac{3x_2}{2} - 4x_3 + \frac{x_6}{2}. \end{cases}$$

The previous operation consisting in maximizing the objective while transforming a non-basic variable into a basic one is called a *pivot*. Moreover x_1 is referred to as the *entering variable* and x_6 as the *leaving variable*.

We now repeat the same process, choosing x_3 which allows the objective to increase faster than x_2 . The third constrain being the tightest one we rewrite it $x_3 = 3/2 - 3x_2/8 - x_5/4 + x_6/8$. The linear program becomes

$$\left\{ \begin{array}{l} z = \frac{111}{4} + \frac{x_2}{16} - \frac{x_5}{8} - \frac{11x_6}{16} \\ x_1 = \frac{33}{4} - \frac{x_2}{16} + \frac{x_5}{8} - \frac{5x_6}{16} \\ x_3 = \frac{3}{2} - \frac{3x_2}{8} - \frac{x_5}{4} + \frac{x_6}{8} \\ x_4 = \frac{69}{4} + \frac{3x_2}{16} + \frac{5x_5}{8} - \frac{x_6}{16} \end{array} \right.$$

The objective can now only be increased by adjusting x_2 , and the tightest constraint being the second we get

$$\left\{ \begin{array}{l} z = 28 - \frac{x_3}{6} - \frac{x_5}{6} - \frac{2x_6}{3} \\ x_1 = 8 + \frac{x_3}{6} + \frac{x_5}{6} - \frac{x_6}{3} \\ x_2 = 4 - \frac{8x_3}{3} - \frac{2x_5}{3} + \frac{x_6}{3} \\ x_4 = 18 - \frac{x_3}{2} + \frac{x_5}{2} \end{array} \right.$$

All the coefficients in the objective function now being negative the solution $(8, 4, 0, 18, 0, 0)$ with objective value 28 is optimal. Therefore, for the original linear problem x_1 , x_2 , and x_3 are 8, 4 and 0, respectively. On their side the slack variables measure how much difference remains in each equality.

Algorithm.

Input : a linear problem in slack form (N, B, A, b, c, v) , l and e the indices of the leaving and entering variables, respectively

Output : an equivalent linear problem in slack form

```

1 Function Pivot( $N, B, A, b, c, v, l, e$ ):
2    $b'_e \leftarrow b_l / a_{le}$ ;
3   foreach  $j \in N \setminus \{e\}$  do  $a'_{ej} \leftarrow a_{lj} / a_{le}$  ;
4    $a'_{el} \leftarrow 1 / a_{le}$ ;
5   foreach  $i \in B \setminus \{l\}$  do
6      $b'_i \leftarrow b_i - a_{ie} b'_e$ ;
7     foreach  $j \in N \setminus \{e\}$  do  $a'_{ij} \leftarrow a_{ij} - a_{ie} a'_{ej}$ ;
8      $a'_{il} \leftarrow -a_{ie} a'_{el}$ 
9   end foreach
10   $v \leftarrow v + c_e b'_e$ ;
11  foreach  $j \in N \setminus \{e\}$  do  $c'_j \leftarrow c_j - c_e a'_{ej}$  ;
12   $c'_l \leftarrow -c_e a'_{el}$ ;
13   $N \leftarrow N \setminus \{e\} \cup \{l\}$ ;  $B \leftarrow B \setminus \{l\} \cup \{e\}$ ;
14  return  $N, B, A', b', c', v$                                /*  $A'$   $m \times n$  matrix of the  $a'_{ij}$  */
15 end
```

When looking at the general process illustrated in the previous example (6.266) a pivot operation only represents one step of a more complex process. In particular all the calls to the `pivot` function are coordinated by the main procedure, the simplex algorithm.

Also note that the problem solved in the example is given in the slack form. Therefore it is also necessary to initialize the algorithm by transforming a problem from standard form into slack form.

We now provide the pseudocode corresponding to the simplex algorithm, and then look at its initialization.

Algorithm. (*Simplex*)

Input : a linear problem in standard form (A, b, c)

Output : an optimal solution to the linear problem

```
1  $(N, B, A, b, c, v) \leftarrow \text{Init}(A, b, c);$ 
2 while  $\exists j \in N$  with  $c_j > 0$  do
3   | select  $e$  such that  $c_e > 0$ ;
4   | foreach index  $i \in B$  do
5   |   | if  $a_{ie} > 0$  then  $d_i \leftarrow b_i/a_{ie}$  else  $d_i \leftarrow \infty$ ;
6   | end foreach
7   | select  $l$  such that  $d_l$  is minimal;
8   | if  $d_l \neq \infty$  then  $(N, B, A, b, c, v) \leftarrow \text{Pivot}(N, B, A, b, c, v, l, e);$ 
9   | else return unbounded;
10 end while
11 for  $i \leftarrow 1$  to  $n$  do
12   | if  $i \in B$  then  $x_i \leftarrow b_i$  else  $x_i \leftarrow 0$ ;
13 end for
14 return  $x_1, \dots, x_n$ 
```

The simplex algorithm (6.272) starts by calling an initialization function. Its role is not only to translate the linear program from standard form into slack form, but also to ensure its feasibility. At times it might be feasible although the basic solution is not, thus requiring more work.

Lemma

Let L be a linear program in standard form. Let x_0 be an additional variable and L_0 be the following linear program with $n + 1$ variables. Maximize $-x_0$, subject to

$$\begin{cases} \sum_{j=1}^n a_{ij}x_j - x_0 \leq b_i, & 1 \leq i \leq m; \\ x_j \geq 0, & 0 \leq j \leq n. \end{cases}$$

Then the linear program L is feasible if and only if the optimal objective value for L_0 is 0.

Example.

Maximize $2x_1 - x_2$, subject to

$$\begin{cases} 2x_1 - x_2 \leq 2 \\ x_1 - 5x_2 \leq -4 \\ x_1, x_2 \geq 0. \end{cases}$$

If we follow the basic strategy to translate the linear program into slack form then this yields the basic solution $x_1 = x_2 = 0$ which violates the second constraint. Therefore we now formulate an *auxiliary linear program* for which we can find a slack form with a feasible basic solution. Applying lemma 6.273 we get the problem:

Maximize $-x_0$, subject to

$$\begin{cases} 2x_1 - x_2 - x_0 \leq 2 \\ x_1 - 5x_2 - x_0 \leq -4 \\ x_1, x_2, x_0 \geq 0. \end{cases}$$

By lemma 6.273 the existence of a feasible solution to the original linear program depends on the optimal value of the objective being 0 or not, for the new linear program. The latter can be rewritten in slack form as

$$\begin{cases} z = & & -x_0 \\ x_3 = 2 - 2x_1 + x_2 + x_0 \\ x_4 = -4 - x_1 + 5x_2 + x_0. \end{cases}$$

This auxiliary linear problem still has no basic solution since it would set x_4 to -4 . However we can reshape it by running the `Pivot` function. The entering variable chosen is x_0 and the leaving one is the basic variable corresponding to the constraint featuring the smallest b_i , i.e. x_4 in our case.

This yields the following slack form

$$\begin{cases} z = -4 - x_1 + 5x_2 - x_4 \\ x_0 = 4 + x_1 - 5x_2 + x_4 \\ x_3 = 6 - x_1 - 4x_2 + x_4. \end{cases}$$

Although the associated basic solution $(4, 0, 0, 6, 0)$ is now feasible, it remains to determine an optimal solution for L_0 . In particular it is sufficient to apply a single call to `Pivot` with entering and leaving variables x_2 and x_0 , respectively. This yields

$$\begin{cases} z = -4 - x_1 + 5x_2 - x_4 \\ x_2 = \frac{4}{5} - \frac{x_0}{5} + \frac{x_1}{5} + \frac{x_4}{5} \\ x_3 = \frac{14}{5} + \frac{4x_0}{5} - \frac{9x_1}{5} + \frac{x_4}{5}. \end{cases}$$

As the objective value for the auxiliary linear program is 0 the initial problem is feasible by lemma 6.273. This also means that x_0 can be removed from the constraints since it is equal to 0.

The last stage consists in restoring the initial objective function ensuring it only contains non-basic variables. We get

$$2x_1 - x_2 = 2x_1 - \left(\frac{4}{5} + \frac{x_1}{5} + \frac{x_4}{5} \right).$$

This yields the following slack form which has a feasible solution

$$\begin{cases} z = -\frac{4}{5} + \frac{9x_1}{5} - \frac{x_4}{5} \\ x_2 = \frac{4}{5} + \frac{x_1}{5} + \frac{x_4}{5} \\ x_3 = \frac{14}{5} - \frac{9x_1}{5} + \frac{x_4}{5}. \end{cases}$$

Algorithm.

Input : a linear problem in standard form (A, b, c)

Output : an initial feasible slack form or infeasible

```

1 Function Init( $A, b, c$ ):
2   select  $k$  such that  $b_k$  is minimal among all the  $b_i$ ;
3   if  $b_k \geq 0$  then return  $\{1, \dots, n\}, \{n+1, \dots, n+m\}, A, b, c, 0$ ;
4   foreach constraint do add  $x_0$  on the right side;
5   set the objective function to  $-x_0$ ;                                     /* define  $L_0$  */
6    $(N, B, A, b, c, v) \leftarrow L_0$  in slack form;
7    $l \leftarrow n + k$ ;
8    $(N, B, A, b, c, v) \leftarrow \text{Pivot}(N, B, A, b, c, v, l, 0)$ ;
9   Run algorithm 6.272 (lines 2-10) until an optimal solution is found for  $L_0$ ;
10  if the optimal solution to  $L_0$  sets  $x_0$  to 0 then
11    if  $x_0$  is basic then perform one pivot to make it non-basic;
12    remove  $x_0$  from the constraints, restore the original objective function, and replace
        any basic variable from the objective function by the right hand side of its
        associated constraint;
13    return the modified final slack form;
14  else
15    return infeasible;
16  end if
17 end

```

Theorem (Fundamental theorem of Linear Programming)

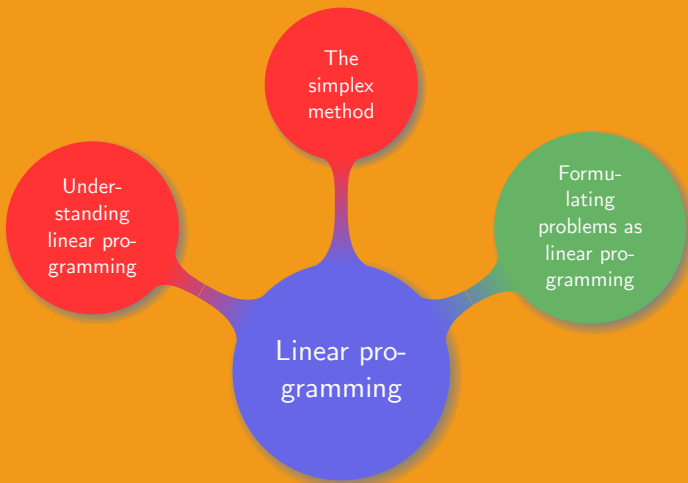
For a linear program L in standard form, one of the following holds.

- L has an optimal solution with a finite objective value;
- L is infeasible;
- L is unbounded;

In any case, the Simplex algorithm returns the appropriate state.

A few words on the Simplex method:

- Exponential time complexity
- Very efficient in practice
- Many variations exist
- For all the variations there exists a family of linear programs which takes exponential time to solve
- The existence of a variation running in polynomial or sub-exponential time remains an open problem



Given two vertices s and t in a weighted graph $G = \langle V, E \rangle$, the shortest weighted path problem (3.149) consists in computing the shortest distance d_t between s and t .

To phrase this problem as a Linear Program it is necessary to determine a set of variables and constraints that define when a shortest path is discovered.

This is in fact what the Bellman-Ford algorithm (3.154) does, as when it terminates the distance from s to any vertex v is given by $d_v \leq d_u + w(u, v)$, where u is a vertex adjacent to v and $w(u, v)$ is the weight of the edge (u, v) (lemma 3.153). Noting that $d_s = 0$ this translates into the following linear program:

Maximize d_t , subject to

$$\begin{cases} d_v \leq d_u + w(u, v) & \text{for each edge } (u, v) \in E \\ d_s = 0 \end{cases} \quad (6.4)$$

The reason behind the maximization of the objective function is again provided by lemma 3.153. The solution to the shortest path problem sets the minimum distance for each vertex v to

$$d_v = \min_{(u,v) \in E} \{d_u + w(u, v)\}.$$

Therefore d_v is the largest value that is less or equal to all the values in $\{d_u + w(u, v)\}$. In other words we are interested in the minimum of the maximum. We want to maximize d_v for all the vertices v on the shortest path from s to t , as it will maximize d_t .

The Linear Program (6.4) can be solved using the simplex algorithm but this results in a slower method than using a specialized algorithm such as the Bellman-Ford (3.154).

Definition

Given a Linear Program where the objective is to be maximized, its *dual* has the same optimal value and is described as follows, using the notations from definition 6.256. Find m real numbers y_1, \dots, y_m that minimize $\sum_{i=1}^m b_i y_i$, while being subject to

$$\begin{cases} \sum_{i=1}^m a_{ij} y_i \geq c_j, & 1 \leq j \leq n; \\ y_i \geq 0, & 1 \leq i \leq m, \end{cases}$$

The original Linear Program is called the *primal*.

Example.

Determine the dual of the following Linear program.

Maximize $3x_1 + x_2 + 2x_3$, subject to

$$\begin{cases} x_1 + x_2 + 3x_3 \leq 30 \\ 2x_1 + 2x_2 + 5x_3 \leq 24 \\ 4x_1 + x_2 + 2x_3 \leq 36 \\ x_1, x_2, x_3 \geq 0 \end{cases}$$

Example.

Determine the dual of the following Linear program.

Maximize $3x_1 + x_2 + 2x_3$, subject to

$$\begin{cases} x_1 + x_2 + 3x_3 \leq 30 \\ 2x_1 + 2x_2 + 5x_3 \leq 24 \\ 4x_1 + x_2 + 2x_3 \leq 36 \\ x_1, x_2, x_3 \geq 0 \end{cases}$$

Minimize $30y_1 + 24y_2 + 36y_3$, subject to

$$\begin{cases} y_1 + 2y_2 + 4y_3 \geq 3 \\ y_1 + 2y_2 + y_3 \geq 1 \\ 3y_1 + 5y_2 + 2y_3 \geq 2 \\ y_1, y_2, y_3 \geq 0. \end{cases}$$

Theorem (Linear Programming duality)

Let (A, b, c) be a primal Linear Program and $x = (x_1, \dots, x_n)$ be a solution returned by the Simplex algorithm. Denoting by c' the vector representing the coefficients in the final slack form, the vector y is defined by $y_i = -c'_{n+i}$ if $n+i$ is in N and zero otherwise. Then x is an optimal solution to the primal Linear Program, y is an optimal solution to the dual Linear Program and

$$\sum_{j=1}^n c_j x_j = \sum_{i=1}^m b_i y_i.$$

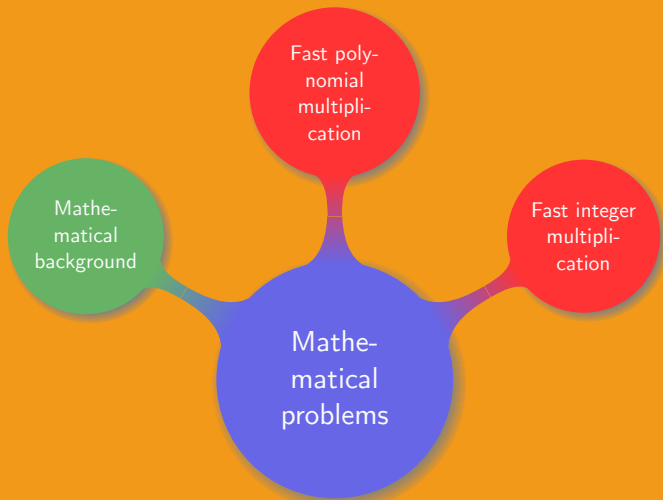
Given a linear program, solve it in:

- Primal form if it has more variables than constraints
- Dual form if it has more constraints than variables

Non-linear optimization function constraints:

- This is not a linear program
- Use a different approach
- If the objective function is quadratic, specific cases can be efficiently solved
- In the general case quadratic programming is \mathcal{NP} -complete

7. Mathematical problems



Many applications require large numbers to be multiplied.

Common multiplication algorithms:

- Simple strategy: $\mathcal{O}(n^2)$
- Karatsuba: $\mathcal{O}(n^{\log_2 3})$

Many applications require large numbers to be multiplied.

Common multiplication algorithms:

- Simple strategy: $\mathcal{O}(n^2)$
- Karatsuba: $\mathcal{O}(n^{\log_2 3})$

Fast Fourier Transform (FFT):

- Fast polynomial and number multiplication
- One of the most important and used algorithms

Definition (Group)

A *group* is a pair (G, \circ) consisting of a set G and a *group operation* $\circ: G \times G \rightarrow G$ that verifies the following properties:

- i *Associativity*: $a \circ (b \circ c) = (a \circ b) \circ c$ for all $a, b, c \in G$
- ii *Existence of a unit element*: there exists an element $e \in G$ such that $a \circ e = e \circ a = a$ for all $a \in G$
- iii *Existence of inverse*: for every $a \in G$ there exists an element $a^{-1} \in G$ such that $a \circ a^{-1} = a^{-1} \circ a = e$

A group is called *abelian* if in addition to the above properties

- iv *Commutativity*: $a \circ b = b \circ a$ for all $a, b \in G$.

Definition (Ring)

A *ring* is a triple $(R, +, \cdot)$ consisting of a set G and two *binary operations* $+, \cdot : R \times R \rightarrow R$ such that

- i $(R, +)$ is an abelian group
- ii *Multiplicative unit*: there exists an element $1 \in G$ such that

$$a \cdot 1 = 1 \cdot a = a \quad \text{for all } a \in R$$

- iii *Associativity*: for any $a, b, c \in R$,

$$a \cdot (b \cdot c) = (a \cdot b) \cdot c$$

- iv *Distributivity*: for any $a, b, c \in R$,

$$a \cdot (b + c) = (a \cdot b) + (a \cdot c), \quad (b + c) \cdot a = (b \cdot a) + (c \cdot a)$$

A ring is called *commutative* if in addition to the above properties

- v *Commutativity*: $a \cdot b = b \cdot a$ for all $a, b \in R$

Definition (Field)

Let $(F, +, \cdot)$ be a commutative ring with unit element of addition 0 and unit element of multiplication 1. Then F is a *field* if

- i $0 \neq 1$
- ii For every $a \in F \setminus \{0\}$ there exists an element a^{-1} such that

$$a \cdot a^{-1} = 1.$$

Definition (Field)

Let $(F, +, \cdot)$ be a commutative ring with unit element of addition 0 and unit element of multiplication 1. Then F is a *field* if

- i $0 \neq 1$
- ii For every $a \in F \setminus \{0\}$ there exists an element a^{-1} such that

$$a \cdot a^{-1} = 1.$$

Remark.

Another way of writing this definition is to say that $(F, +, \cdot)$ is a field if $(F, +)$ and $(F \setminus \{0\}, \cdot)$ are abelian groups, \cdot distributes over $+$, and $0 \neq 1$.

Example.

Let n be an integer, and $\mathbb{Z}/n\mathbb{Z}$ be the set of the integers modulo n

- $(\mathbb{Z}/n\mathbb{Z}, +)$ also denoted $(\mathbb{Z}_n, +)$ is a group
- $(\mathbb{Z}/n\mathbb{Z}, +, \cdot)$ is a ring
- If n is prime then $(\mathbb{Z}/n\mathbb{Z}, +, \cdot)$ is the field \mathbb{F}_n
- The invertible elements of $\mathbb{Z}/n\mathbb{Z}$, with respect to ' \cdot ', form a group denoted $U(\mathbb{Z}/n\mathbb{Z})$ or sometimes \mathbb{Z}_n^\times or \mathbb{Z}_n^*
- $(\mathbb{Z}/n\mathbb{Z}[X], +, \cdot)$ is the ring of the polynomials over $\mathbb{Z}/n\mathbb{Z}$
- If n is prime and the polynomial $P(X)$ is irreducible then $(\mathbb{F}_n[X]/\langle P(X) \rangle, +, \cdot)$ is a field; this is $\mathbb{F}_{n^{\deg P(X)}}$

Definitions

Let R be a ring and n be a strictly positive integer.

- ① An element a of R is a *zero divisor* if there exists $b \in R \setminus \{0\}$ such that $ab = 0$.
- ② If 0 is the only zero divisor in R , then R is an *integral domain*.
- ③ An element $\omega \in R$ is an *n th root of unity* if $\omega^n = 1$.
- ④ An element $\omega \in R$ is a *primitive n th root of unity* if $\omega^n = 1$ and for any $1 \leq k \leq n-1$, $\omega^k \neq 1$.

Example.

- In \mathbb{C} , $e^{2i\pi/8}$ is a primitive 8th root of unity;
- In \mathbb{Z}_{17} , 2 is not a primitive 16th root of unity;

Lemma

Let R be a ring, l, n be two integers such that $1 < l < n$, and ω be a primitive n th root of unity. Then (i) $\omega^l - 1$ is not a zero divisor in R , and (ii) $\sum_{j=0}^{n-1} \omega^{lj} = 0$.

Proof.

(i) By definition of a primitive n th root of unity ω^l is not 1 unless l is 0 or larger or equal to n .

(ii) Note that for any $c \in R$ and m in \mathbb{N}

$$c^m - 1 = (c - 1)(1 + c + c^2 + \cdots + c^{m-1}).$$

In particular for $c = \omega^l$ and $m = n$

$$\omega^{ln} - 1 = (\omega^l - 1)(1 + \omega^l + \cdots + \omega^{l(n-1)}).$$

As $\omega^{ln} = 1$, and $\omega^l - 1$ is not a zero divisor, $\sum_{j=0}^{n-1} \omega^{lj} = 0$. □

Definition

Let R be a ring, and $\omega \in R$ be a primitive n th root of unity. We denote a polynomial $P(X)$ of degree less than n in $R[X]$ by its coefficients

$$P(X) = \sum_{i=0}^{n-1} a_i X^i = (a_0, \dots, a_{n-1}).$$

The linear map

$$\begin{aligned} \text{DFT}_\omega : R^n &\longrightarrow R^n \\ (a_0, \dots, a_{n-1}) &\longmapsto (P(1), P(\omega), \dots, P(\omega^{n-1})) \end{aligned}$$

evaluates P at the powers of ω and is called *Discrete Fourier Transform* (DFT).

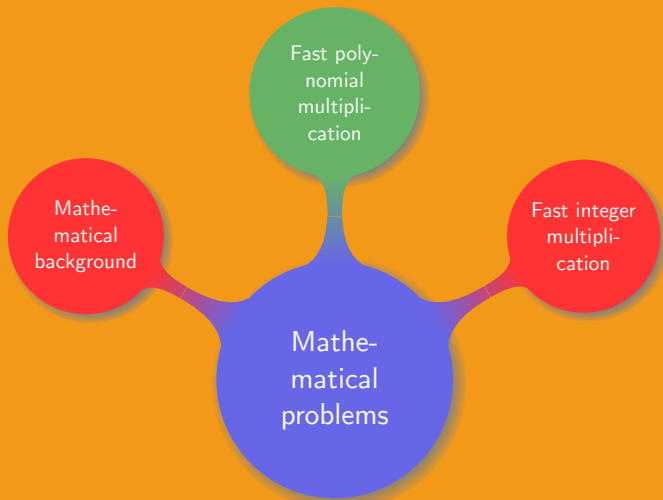
As DFT_ω is a linear map it is expressed as a matrix transformation

$$\begin{pmatrix} P(1) \\ P(\omega) \\ P(\omega^2) \\ \vdots \\ P(\omega^{n-1}) \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega & \omega^2 & \cdots & \omega^{n-1} \\ 1 & \omega^2 & \omega^4 & \cdots & \omega^{2(n-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{n-1} & \omega^{2(n-1)} & \cdots & \omega^{(n-1)^2} \end{pmatrix}}_{V_\omega} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{pmatrix}.$$

Note that if ω is a primitive n th root of unity then so is ω^{-1} . Otherwise there would be a $1 \leq k < n$ such that $(\omega^{-1})^k$ is 1. And as $\omega^n = 1$ we would have $1 = \omega^{n+k}$ and $\omega^k = 1$.

Then using lemma 7.297 observe that $V_\omega V_{\omega^{-1}} = nI_n$, where I_n is the identity matrix of size $n \times n$. Thus the inverse of DFT_ω is

$$\text{DFT}_\omega^{-1} = \frac{1}{n} \text{DFT}_{\omega^{-1}}.$$



Two main cases depending on the structure of the polynomial:

- Dense: use an array where the coefficient of each monomial is stored at index “degree of the monomial”
- Sparse: use a structure composed of two arrays storing the degrees and the corresponding coefficients, respectively

Two main cases depending on the structure of the polynomial:

- Dense: use an array where the coefficient of each monomial is stored at index “degree of the monomial”
- Sparse: use a structure composed of two arrays storing the degrees and the corresponding coefficients, respectively

Alternative strategy: over an integral domain a polynomial of degree strictly less than n can be represented using its value at n distinct points

Let $P(X) = \sum_{i=0}^n a_i X^i$ be a polynomial of degree n . Evaluating P costs $\mathcal{O}(n^2)$ if naively computed. However note that $P(X)$ can be rewritten

$$P(X) = a_0 + X(a_1 + X(a_2 + \cdots + X(a_{n-1} + Xa_n))).$$

This remark dramatically decreases the complexity as it drops to $\mathcal{O}(n)$, and yields the following simple algorithm.

Algorithm. (*Horner*)

Input : a polynomial P and x the value to evaluate P at

Output: Px the evaluation of P at x

```
1 Function Horner( $P, x$ ):  
2    $Px \leftarrow 0$ ;  
3   for  $i \leftarrow \deg P$  to 0 do  $Px \leftarrow Px \cdot x + \text{coeff}[i]$  ;  
4   return  $Px$ ;  
5 end
```

Two main cases depending on the polynomial representation:

- Dense: usual approach, multiply the various coefficients together; complexity $\mathcal{O}(n^2)$
- Evaluation: evaluate the polynomials in n points, multiply them two by two as

$$PQ(x_i) = P(x_i)Q(x_i).$$

Complexity is $\Omega(n^2)$ since n evaluations are necessary, to which have to be added the cost of the multiplications and of the interpolation to get the final polynomial (usually $\mathcal{O}(n^2)$).

Looking back at DFT_ω , it can be viewed as a special multipoint evaluation of a polynomial in the powers $1, \omega, \dots, \omega^{n-1}$ of a primitive n th root of unity ω . Then its inverse DFT_ω^{-1} , which given n evaluations of a polynomial allows to recover its coefficients, is just the interpolation at the powers of ω .

From the previous discussion on the DFT (7.299), it is clear that knowing how to compute it efficiently means being able to also compute its inverse efficiently.

For the sake of simplicity assume $n = 2^k$, $k \in \mathbb{N}$, and observe that

$$\begin{aligned} P(X) &= \sum_{i=0}^{n-1} a_i X^i \\ &= (a_0 + a_2 X^2 + \dots + a_{n-2} X^{n-2}) + (a_1 X + a_3 X^3 + \dots + a_{n-1} X^{n-1}) \\ &= P_1(X^2) + X P_2(X^2) \end{aligned} \tag{7.1}$$

with both P_1 and P_2 of degree less than $(n-2)/2 < n/2$.

The structure of equation (7.1) suggests a “divide and conquer” approach in order to determine

$$P(\omega^i) = P_1(\omega^{2i}) + \omega^i P_2(\omega^{2i}), \quad 0 \leq i < n. \quad (7.2)$$

This formulation can be further rewritten by noticing that

$$\begin{aligned} 0 &= \omega^n - 1 \\ &= (\omega^{n/2} - 1)(\omega^{n/2} + 1). \end{aligned}$$

By lemma 7.297 none of the two factors is a zero divisor and as $\omega^{n/2} \neq 1$, ω being a primitive n th root of the unity, we conclude that $\omega^{n/2} = -1$. Hence, for all $0 \leq i < n/2$, $\omega^i = -\omega^{n/2+i}$, and (7.2) can be rewritten

$$\begin{aligned} P(\omega^i) &= P_1(\omega^{2i}) + \omega^i P_2(\omega^{2i}), \quad 0 \leq i < n/2, \\ P(\omega^{n/2+i}) &= P_1(\omega^{2i}) - \omega^i P_2(\omega^{2i}), \quad 0 \leq i < n/2. \end{aligned} \quad (7.3)$$

Algorithm. (*Fast Fourier Transform (FFT)*)

Input : a polynomial P of degree $< n$, with n a power of 2, and ω a primitive n th root of unity

Output : $\text{DFT}_\omega(P)$

```
1 Function FFT( $P, \omega$ ):  
2    $n \leftarrow \deg P + 1$ ;  
3   if  $n = 1$  then return  $P$ ;  
4    $P_1 \leftarrow \text{FFT}(\sum_{j=0}^{n/2-1} a_{2j} X^{2j}, \omega^2)$ ;  
5    $P_2 \leftarrow \text{FFT}(\sum_{j=0}^{n/2-1} a_{2j+1} \omega X^{2j}, \omega^2)$ ;  
6   for  $i \leftarrow 0$  to  $n/2$  do  
7      $P_\omega[i] \leftarrow P_{1,\omega^2}[i] + \omega^i P_{2,\omega^2}[i]$ ;  
8      $P_\omega[n/2 + i] \leftarrow P_{1,\omega^2}[i] - \omega^i P_{2,\omega^2}[i]$ ;  
9   end for  
10  return  $P_\omega$ ;  
11 end
```

Theorem

Given a polynomial P over a commutative ring and ω a primitive n th root of unity, the FFT algorithm correctly computes $\text{DFT}_\omega(P)$ in time $\mathcal{O}(n \log n)$.

Proof.

The correctness clearly follows from the previous discussion, more particularly from the recurrence relation (7.3).

Let $T(n)$ be the time to compute a DFT. Then in relation (7.3) two smaller DFT are computed, plus $n/2$ multiplications and n additions, that is

$$\begin{aligned} T(n) &\leq 2T(n/2) + n/2 + n \\ &\leq 2T(n/2) + 3n/2. \end{aligned}$$

By the Master theorem (2.95) the time complexity of a DFT is $\mathcal{O}(n \log n)$. □

Let R be a ring containing a primitive n th root of unity. Given two polynomials P and Q with degrees less than $n/2$, defined over $R[X]$, we want to determine $S = PQ$. Note that n is still taken to be a power of 2.

As both P and Q are of degrees less than n they can be efficiently evaluated using the FFT algorithm (7.306). Then n multiplications in R are enough to determine S , represented using its evaluation in n points.

Applying DFT_{ω}^{-1} to the n evaluations of S , is achieved through the calculation of $1/n \text{DFT}_{\omega^{-1}} S$ (7.299). This computation returns the interpolation of S in n points (7.304), that is it determines the unique polynomial of degree less than n passing through the n points. Hence we obtain a fast strategy to compute the product of two polynomials.

Algorithm. (*Fast polynomial multiplication*)

Input : P and Q two polynomials of degree $< n/2$, with n a power of 2,
a primitive n th root of unity ω

Output : $S = PQ$

```
1 Function FPMult( $P, Q, \omega$ ):  
2    $P_\omega \leftarrow \text{FFT}(P, \omega);$   
3    $Q_\omega \leftarrow \text{FFT}(Q, \omega);$   
4    $S_\omega \leftarrow P_\omega Q_\omega;$   
5    $S \leftarrow \frac{1}{n} \text{FFT}(S_\omega, \omega^{-1});$   
6   return  $S$   
7 end
```

Algorithm. (*Fast polynomial multiplication*)

Input : P and Q two polynomials of degree $< n/2$, with n a power of 2,
a primitive n th root of unity ω

Output : $S = PQ$

```
1 Function FPMult( $P, Q, \omega$ ):  
2    $P_\omega \leftarrow \text{FFT}(P, \omega);$   
3    $Q_\omega \leftarrow \text{FFT}(Q, \omega);$   
4    $S_\omega \leftarrow P_\omega Q_\omega;$   
5    $S \leftarrow \frac{1}{n} \text{FFT}(S_\omega, \omega^{-1});$   
6   return  $S$   
7 end
```

Definition

A commutative ring containing a primitive 2^k th root of unity for any k in \mathbb{N}^* is said to *support the FFT*.

Theorem

Let R be a ring supporting the FFT and n be 2^k , with k in \mathbb{N}^* . Then for two polynomials P and Q in $R[X]$, with $\deg PQ < n$, the fast polynomial multiplication algorithm computes their product in time $\mathcal{O}(n \log n)$.

Proof.

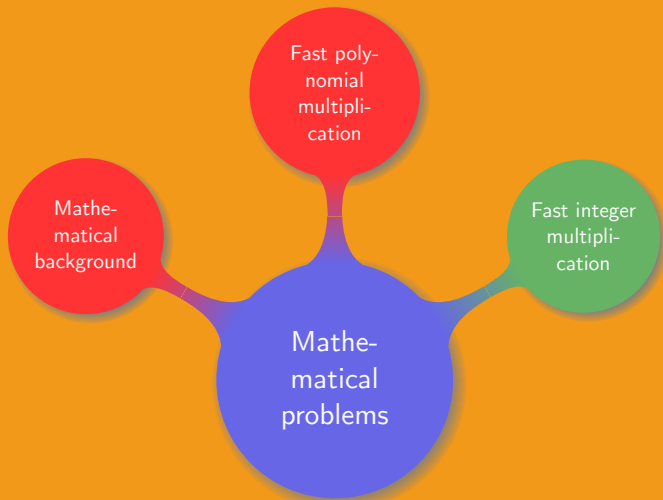
The correctness of the algorithm is clear when considering the previous discussion (7.308).

The algorithm computes three DFT, n component-wise products in R , as well as n multiplications by the inverse of n in R . Therefore the overall complexity is dominated by $\mathcal{O}(n \log n)$.



In order to run the fast polynomial multiplication algorithm (7.309) the underlying ring R must support the FFT. In the case where R does not contain any primitive 2^k th root of unity a “virtual” one can be attached to the ring.

The Schönage-Strassen algorithm handles this special case at the cost of a slightly worse complexity. In fact their result states that over any ring the product of two polynomials of degree less than n can be computed in $\mathcal{O}(n \log n \log \log n)$ operations.



Let a and b be two N -bit long integers. For the sake of simplicity we assume N to be of the form 2^{2^l} for some integer $l > 0$. Let a_N, \dots, a_0 and b_N, \dots, b_0 denote the binary representations of a and b , respectively.

Since N was chosen to be a square it is easy to split a and b into blocks of size \sqrt{N} and write them

$$a = \sum_{i=0}^{\sqrt{N}-1} A_i 2^{i\sqrt{N}}, \text{ and } b = \sum_{i=0}^{\sqrt{N}-1} B_i 2^{i\sqrt{N}}, \quad 0 \leq A_i, B_i \leq 2^{\sqrt{N}} - 1.$$

If we consider the polynomials $A(X) = \sum_{i=0}^{\sqrt{N}-1} A_i X^i$ and $B(X) = \sum_{i=0}^{\sqrt{N}-1} B_i X^i$, then the product AB evaluated at $2^{\sqrt{N}}$ is exactly the product ab . Therefore integer multiplication can be performed via polynomial multiplication: (i) write the two integers as polynomials, (ii) apply the fast polynomial multiplication on them, and (iii) finally evaluate the product at $2^{\sqrt{N}}$.

While the first step is simple to achieve the second one requires more technical considerations. In fact the two polynomials A and B are both of degree less than $2\sqrt{N}$ and as such R must contain a primitive $2\sqrt{N}$ th root of unity.

As \mathbb{Z} does not support FFT some extra work is needed in order to attach a new “virtual” element to the ring without altering the final result.

Consider the ring $\mathbb{Z}_{2^{\sqrt{N}+1}}$ and observe that 2 is a primitive $2\sqrt{N}$ th root of unity. This is clear as $2^{\sqrt{N}}$ is -1 modulo $2^{\sqrt{N}+1}$. Thus $t = 2\sqrt{N}$ is the smallest power for which 2^t is 1.

An obvious idea is then to perform the computation in the ring $\mathbb{Z}_{2^{\sqrt{N}+1}}[X]$. However as both A and B can feature coefficients as large as $2^{\sqrt{N}} - 1$, the polynomial C resulting from their product can have coefficients up to $\sqrt{N}2^{2\sqrt{N}}$, which is larger than $2^{\sqrt{N}+1}$.

As a result, if coefficients in C happen to be too large they will be reduced modulo $2^{\sqrt{N}}+1$, ruining the whole calculation. A simple solution consists in performing all the computation in a larger ring. At that stage two points must be taken into consideration: (i) the use of a large ring increases the computational cost and (ii) the ring must contain a primitive $2\sqrt{N}$ th root of unity.

Note that for any $N \geq 1$, $2^{3\sqrt{N}} > \sqrt{N}2^{2\sqrt{N}}$, while 8 is a primitive $2\sqrt{N}$ th root of unity in $\mathbb{Z}_{2^{3\sqrt{N}}+1}$. The latter being a consequence of 2 being a primitive $6\sqrt{N}$ th root of unity.

Therefore it suffices to consider A and B as polynomials over the ring $\mathbb{Z}_{2^{3\sqrt{N}}+1}[X]$. Then applying the fast polynomial multiplication algorithm (7.309) and evaluating the product at $2^{\sqrt{N}}$ yields the result.

Algorithm. (*Fast integer multiplication*)

Input : two N bit integers a and b , a ring $R = \mathbb{Z}_{2^{3\sqrt{N}}+1}$, $\omega = 8$ a primitive $2\sqrt{N}$ th root of unity

Output: $c = a \cdot b$

```
1  $A \leftarrow \text{poly}(a);$                                 /* encode  $a$  as a polynomial */
2  $B \leftarrow \text{poly}(b);$                                 /* encode  $b$  as a polynomial */
3  $C \leftarrow \text{FPMult}(A, B, \omega);$ 
4  $c \leftarrow \text{Horner}(C, 2^{\sqrt{N}});$ 
5 return  $c;$ 
```

Theorem

Given two integers of length N in a ring R , the fast integer multiplication algorithm computes their product in $\mathcal{O}(\sqrt{N} \log \sqrt{N})$ arithmetic operations in R .

Proof.

The correctness results from the previous discussion (7.313).

The pre-dominant computation in the algorithm is the fast polynomial multiplication of A and B , which by theorem 7.310 takes $\mathcal{O}(\sqrt{N} \log \sqrt{N})$. \square

Remark.

With a bit more work it is possible to determine the complexity in term of bit operations, instead of arithmetic operations. In that case the complexity becomes

$$\mathcal{O}(N \log^{2+\log_2 3-1} N).$$

Other interesting remarks:

- At the bit level, most operations can be performed using “shifts”, incurring a linear cost in the length of the integers. This is, in particular the main reason for choosing ω to be a power of 2.
- The integer $2\sqrt{N}$ has inverse $2^{6\sqrt{N}-\log_2 \sqrt{N}-1}$ in $\mathbb{Z}_{2^{3\sqrt{N}+1}}$. Observe that $2^{6\sqrt{N}} \equiv 1 \pmod{2^{3\sqrt{N}+1}}$, and $2^{\log_2 \sqrt{N}} = \sqrt{N}$.
- Although in the algorithm no coefficient reaches $2^{3\sqrt{N}} + 1$ all the calculations are performed in the ring $\mathbb{Z}_{2^{3\sqrt{N}+1}}$ and the special case of adding two $3\sqrt{N}$ bits long elements has to be considered when defining addition.
- To date the asymptotically fastest integer multiplication algorithm, due to Fürer, takes $N \log N 2^{\mathcal{O}(\log^* N)}$ bit operations.

Definition

The complexity of multiplying two polynomials of degree less than n is denoted $M(n)$, while the complexity of multiplying two n -bit integer is denoted $M_I(n)$.

For all n and m in \mathbb{N} ,

$$M(n + m) \geq M(m) + M(n) \quad \text{and} \quad M_I(n + m) \geq M_I(n) + M_I(m).$$

These simpler notations allows an easier complexity study of more advanced algorithms such as fast multi-point evaluation and fast interpolation.

8. More advanced topics



Given an \mathcal{NP} -complete problem:

- It is impossible to efficiently find a solution
- In practice it might have many applications

Hope for such problems:

- Inputs are always small
- Specific sub-cases arising in practice can be solved efficiently
- An almost optimal solution is sufficient and can be computed efficiently

Definition

- ① For a given problem P , an algorithm returning a near-optimal solution to P is called an *approximation algorithm*.
- ② Let the cost of an optimal solution be C^* and the one of an approximation be C . If for any input of size n there exists an *approximation ratio* $\rho(n)$, such that

$$\max\left(\frac{C}{C^*}, \frac{C^*}{C}\right) \leq \rho(n),$$

then the algorithm is said to be a $\rho(n)$ -*approximation algorithm*.

Remark.

- The approximation ratio can never be less than 1
- Approximation algorithms are expected to
 - Be polynomial time
 - Feature slowly growing approximation ratio as n increases
- Some approximation algorithms take an input parameter defining the precision of the approximation. The more precise the approximation, the longer the running time.

Problem (Optimal Vertex Cover)

Let $G = \langle V, E \rangle$ be an undirected graph. A *vertex cover* is a subset $V' \subseteq V$ such that if $(u, v) \in E$ then at least u or v is in V' . Find a vertex cover of minimum size.

Problem (Optimal Vertex Cover)

Let $G = \langle V, E \rangle$ be an undirected graph. A *vertex cover* is a subset $V' \subseteq V$ such that if $(u, v) \in E$ then at least u or v is in V' . Find a vertex cover of minimum size.

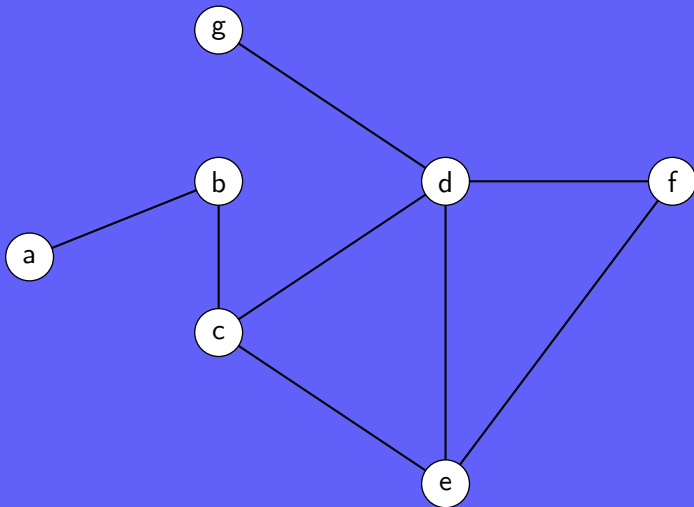
Algorithm. (*Approx Vertex Cover*)

Input : $G = \langle V, E \rangle$ an undirected graph

Output : C a nearly optimal vertex cover

```
1  $C \leftarrow \emptyset$ ;  $E' \leftarrow G.E$ ;  
2 while  $E' \neq \emptyset$  do  
3   |  $e \leftarrow$  an arbitrary edge of  $E'$ ;  $C \leftarrow C \cup \{e.u, e.v\}$ ;  
4   | remove from  $E'$  any edge incident to  $u$  or  $v$ ;  
5 end while  
6 return  $C$ ;
```

Exercise.



Theorem

Approx Vertex Cover runs in time $\mathcal{O}(|V| + |E|)$ and is a 2-approximation algorithm.

Proof.

Representing E' using an adjacency list leads to $\mathcal{O}(|V| + |E|)$.

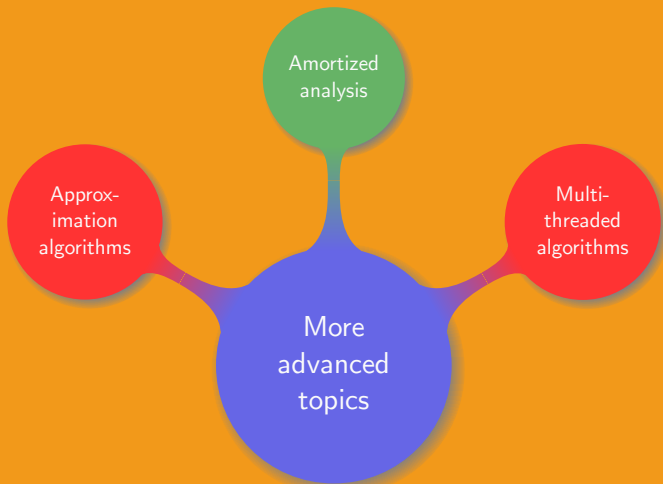
First note that C is a vertex cover since the algorithm loops until all the edges from E have been covered by some vertex in C .

Let A denote the set of all the edges selected by the algorithm. Then any cover includes at least one endpoint, and by construction no two edges in A share a common endpoint. Thus there is no two edges in A that are covered by the same vertex in C^* , and we have the lower bound $|C^*| \geq |A|$.

Finally noting that $|C| = 2|A|$ yields $|C| \leq 2|C^*|$. □

Common methodology for approximation algorithms:

- An optimal solution is unknown
e.g. the size of the vertex cover is unknown
- Determine a lower bound on an optimal solution
e.g. in approx vertex cover (8.326) the set of the selected edges forms a maximum matching (4.210), which provides a lower bound on the size of an optimal vertex cover
- Relate the size of the approximation to the lower bound on the optimal solution
e.g. the approximation ratio is obtained by relating $|C|$ to $|A|$



A *dynamic table* is an array which automatically resizes as elements are added or removed. A common strategy consists in doubling the size of the table as soon as an overflow occurs. In such a context we want to determine the cost of n insertions.

We define the cost c_i of the i th operation to be i , if $i - 1$ is a power of 2, and 1 otherwise. The cost C_n of n insertions is given by

$$\begin{aligned} C_n &= \sum_{i=1}^n c_i \\ &\leq n + \sum_{j=0}^{\lfloor \log(n-1) \rfloor} 2^j \\ &\leq 3n = \Theta(n) \end{aligned}$$

Hence the average cost for each operation is $\Theta(n)/n = \Theta(1)$.

Definition (Amortized analysis)

Given a sequence of operations, an *amortized analysis* is a strategy allowing to show that the average cost per operation is small although single operations in the sequence might be expensive.

Remark.

Amortize analysis:

- Does not use probabilities
- Evaluates the average performance of each operation in the worst case
- Provides a more precise and useful evaluation of the difficulty of a problem than average case complexity

Three main approaches to amortized analysis:

- *Aggregate method*: most simple approach where the total running time for the sequence is analysed and divided by the number of operations
- *Accounting method*: charge each type of operation a constant cost such that the extra charge on inexpensive ones can be stored in a “bank” and used to pay expensive subsequent operations
- *Potential method*: evaluate and store the total amount of extra work done over the whole data structure and release it to cover the cost of subsequent operations

Remark.

The average cost for each operation in the case of dynamic tables (8.331) was determined using the aggregate method.

Let c_i denote the actual real cost of the i th operation, while \hat{c}_i represents its amortized cost, i being larger than 1. The goal is to always ensure that at any stage the sum of the \hat{c}_i is larger than the sum of the c_i , meaning that some extra credit is available for future operations.

For a sequence of n operations the bank balance never becoming negative can be expressed as

$$\sum_{i=1}^n \hat{c}_i - \sum_{i=1}^n c_i \geq 0.$$

Example.

In the case of dynamic tables (8.331) charge $\hat{c}_i = 3$ units for each insertion: one is used on insertion while the 2 remaining are saved for a future use. In particular when memory is reallocated a unit is use for reassigning the current element while the last one is spent to move an older value.

Let D_0 be an initial data structure. At step i , the i th operation, costing c_i , is applied to D_{i-1} . The *potential* associated with the data structure D_i , denoted $\Phi(D_i)$, is a real number and Φ is called the *potential function*. The amortized cost corresponds to the actual cost plus the change in potential induced by the operation. This formalizes as $\hat{c}_i = c_i + \Phi(D_i) - \Phi(D_{i-1})$, and we get

$$\begin{aligned}\sum_{i=1}^n \hat{c}_i &= \sum_{i=1}^n c_i + \Phi(D_i) - \Phi(D_{i-1}) \\ &= \sum_{i=1}^n c_i + \Phi(D_n) - \Phi(D_0) \geq \sum_{i=1}^n c_i\end{aligned}$$

Example.

For the dynamic tables (8.331) set the potential after the i th insertion to $\Phi(D_i) = 2i - 2^{\lceil \log_2 i \rceil}$. Then whether or not $i - 1$ is a power of 2, the amortized cost is 3.

Remarks on the three methods:

- The accounting and potential methods are more powerful and refined than the aggregate one; they are equivalent in terms of applicability and precision of the bound provided
- The accounting method charges a different cost for each type of operation
- The potential method focuses on the effect of a particular operation at a specific time, and in particular on the cost of future operations
- Different methods might lead to different bounds, but always upper bound the actual cost

For the Union-Find data structure (1.56), theorem 2.91 can be proven using either the accounting or the potential method.

The potential method suits this example since the goal is to determine how fast the tree is flattening, or in other words how each operation affects the whole data structure. By observing the change in potential after each type of operation it is then possible to determine the amortized cost associated with each of them.

The first and less straight forward step consists in properly setting the potential function. Two cases have to be considered (i) x is a root or has rank 0, and (ii) x is not a root and has rank larger or equal to 1. Defining the potential of x after i operations is most complicated in the latter case.

Once the potential function has been properly defined it can be bounded by $0 \leq \phi_i(x) \leq \alpha(n) \cdot x.rank$, and the amortized cost of the various operations can be determined.



Algorithm discussed so far targeted *uniprocessor* computers, while most modern systems feature *multiprocessors* sharing a common memory. The question is then to know how to adapt those algorithms to this new context, and in particular how to efficiently partition the work among several *threads*, each having roughly the same load.

The most simple strategy consists in employing a software layer, called *concurrency platform*, which coordinates, schedules, and manages the resources.

A simple extension to *serial programming* is the addition of the following instructions: `parallel`, `spawn`, and `sync`.

The three new instructions:

- `parallel`: added to loops to indicate that iterations can be computed in parallel
- `spawn`: executes a new parallel process, the current one may then choose to run concurrently or wait for its child
- `sync`: requests the process to wait for all spawned processes to complete

Sieve of Eratosthenes

	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33
34	35	36	37	38	39	40	41	42	43	44
45	46	47	48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63	64	65	66
67	68	69	70	71	72	73	74	75	76	77
78	79	80	81	82	83	84	85	86	87	88
89	90	91	92	93	94	95	96	97	98	99
100	101	102	103	104	105	106	107	108	109	110
111	112	113	114	115	116	117	118	119	120	121

Sieve of Eratosthenes

	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33
34	35	36	37	38	39	40	41	42	43	44
45	46	47	48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63	64	65	66
67	68	69	70	71	72	73	74	75	76	77
78	79	80	81	82	83	84	85	86	87	88
89	90	91	92	93	94	95	96	97	98	99
100	101	102	103	104	105	106	107	108	109	110
111	112	113	114	115	116	117	118	119	120	121

Sieve of Eratosthenes

	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33
34	35	36	37	38	39	40	41	42	43	44
45	46	47	48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63	64	65	66
67	68	69	70	71	72	73	74	75	76	77
78	79	80	81	82	83	84	85	86	87	88
89	90	91	92	93	94	95	96	97	98	99
100	101	102	103	104	105	106	107	108	109	110
111	112	113	114	115	116	117	118	119	120	121

Sieve of Eratosthenes

	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33
34	35	36	37	38	39	40	41	42	43	44
45	46	47	48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63	64	65	66
67	68	69	70	71	72	73	74	75	76	77
78	79	80	81	82	83	84	85	86	87	88
89	90	91	92	93	94	95	96	97	98	99
100	101	102	103	104	105	106	107	108	109	110
111	112	113	114	115	116	117	118	119	120	121

Sieve of Eratosthenes

	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33
34	35	36	37	38	39	40	41	42	43	44
45	46	47	48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63	64	65	66
67	68	69	70	71	72	73	74	75	76	77
78	79	80	81	82	83	84	85	86	87	88
89	90	91	92	93	94	95	96	97	98	99
100	101	102	103	104	105	106	107	108	109	110
111	112	113	114	115	116	117	118	119	120	121

Sieve of Eratosthenes

	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31	32	33
34	35	36	37	38	39	40	41	42	43	44
45	46	47	48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63	64	65	66
67	68	69	70	71	72	73	74	75	76	77
78	79	80	81	82	83	84	85	86	87	88
89	90	91	92	93	94	95	96	97	98	99
100	101	102	103	104	105	106	107	108	109	110
111	112	113	114	115	116	117	118	119	120	121

In Eratosthenes sieve all the loop iterations are independent from each others. It is therefore extremely simple to parallelize the algorithm only using the `parallel` keyword in conjunction with the `for` loop.

In other cases a *race condition* might occurs, leading to a result that is non-deterministic. A simple example is as follows.

Algorithm.

Input : x

Output: $x + 2$ is expected

```
1 parallel for  $i \leftarrow 1$  to 2 do  
2   |  $x \leftarrow x + 1$ ;  
3 end for  
4 return  $x$ ;
```

Two main strategies can be employed to avoid race conditions: (i) a thread sets a *lock* when reaching a critical part of the code to prevent any other thread to run it at the same time; (ii) use special hardware instructions called *atomic operations* which can run several operations at once.

Definition (Computation DAG)

Multi-threaded algorithms can be represented as Directed Acyclic Graph (DAG), often referred to as *computation dag*.

More specifically each vertex stands for an instruction while the edges organize the dependencies between the various instructions. An edge (u, v) implies that instruction u must be run before v .

A chain composed of one or more instructions and not containing any concurrency instruction is called a *strand*. Two strands connected by a directed path are in *series*, or otherwise in *parallel*.

A computer composed of n processors is expected to run n concurrent threads. The *scheduler* is the part of the Operating System which decides which thread to run and on which CPU.

Several approaches are to be used depending on the hardware. Modern desktop and laptop computers feature a memory shared among several CPU, which is much easier to handle than having an independent memory for each CPU.

Example.

OpenMP allows to easily add parallelism to existing source code without requiring any significant rewrite. It perfectly suites systems where the memory is shared among multiple CPUs.

MPI offers advanced possibilities more specifically targeting systems with a distributed memory. It is very common in the realm of high end computing, especially on clusters. MPI often requires a complete redesign, or even change of algorithm to implement.

Assuming the existence of a Merge function, the following algorithm performs serial Merge Sort on a given list L .

Algorithm. (*Merge Sort*)

Input : $L = a_1, \dots, a_n$

Output : L , sorted into elements in non-decreasing order

```
1 Function MergeSort( $L$ ):  
2   if  $n > 1$  then  
3      $m \leftarrow \lfloor n/2 \rfloor$ ;  
4      $L_1 \leftarrow a_1, \dots, a_m$ ;  $L_2 \leftarrow a_{m+1}, \dots, a_n$ ;  
5      $L \leftarrow \text{Merge}(\text{MergeSort}(L_1), \text{MergeSort}(L_2))$   
6   end if  
7   return  $L$   
8 end
```

The recursive calls MergeSort seem to make this algorithm a good candidate for parallelism.

Algorithm. (*Merge*)

Input : L_1, L_2 , two sorted lists

Output : L a merged list of L_1 and L_2 , with elements in increasing order

```
1 Function Merge( $L_1, L_2$ ):  
2    $L \leftarrow \emptyset$ ;  
3   while  $L_1 \neq \emptyset$  and  $L_2 \neq \emptyset$  do  
4      $x \leftarrow \min(L_1, L_2)$ ;          /* smallest element in  $L_1$  and  $L_2$  */  
5     append  $x$  to  $L$ ;  
6     if  $L_1 = \emptyset$  or  $L_2 = \emptyset$  then  
7        $L_1 = \emptyset$  ? append  $L_2$  to  $L$  : append  $L_1$  to  $L$ ;  
8     end if  
9   end while  
10  return  $L$ ;  
11 end
```

Algorithm. (*Parallel Merge Sort*)

Input : $L = a_1, \dots, a_n$

Output : L , sorted into elements in non-decreasing order

```
1 Function P-MergeSort( $L$ ):  
2   if  $n > 1$  then  
3      $m \leftarrow \lfloor n/2 \rfloor$ ;  
4      $L_1 \leftarrow a_1, \dots, a_m$ ;  $L_2 \leftarrow a_{m+1}, \dots, a_n$ ;  
5      $L_1 \leftarrow$ spawn P-MergeSort ( $L_1$ );  
6      $L_2 \leftarrow$ spawn P-MergeSort ( $L_2$ );  
7     sync  $L \leftarrow$ Merge ( $L_1, L_2$ );  
8     return  $L$ ;  
9   end if  
10 end
```

In fact this parallel version does not improve much on the serial version: although Merge Sort is parallel, Merge remains serial and as such the recurrence relation still features the same $\mathcal{O}(n)$ term. Moreover the parallelism is only $\Theta(\log n)$.

Although at times serial algorithms can easily be made parallel, it is often impossible, and they require to be completely redesigned using a totally different approach. In some other cases algorithms might seem to be intrinsically serial, but can still be adjusted at the cost of some extra work. This is for instance the case of the Merge algorithm.

The idea is to carefully generate four lists that can be merged two by two without altering the ordering. The process is as follows.

- 1 Find the median element in the longest list: define two sublists $L_{1,L}$ and $L_{1,R}$
- 2 Determine its corresponding potential location in the second list: define two sublists $L_{2,L}$ and $L_{2,R}$
- 3 Recursively merge the lists $L_{1,L}$ and $L_{2,L}$ and $L_{1,R}$ and $L_{2,R}$

Algorithm. (*Parallel Merge*)

Input : L_1, L_2 , two sorted lists, with $L_1.length > L_2.length$

Output : L , a merged list of L_1 and L_2 , with elements in increasing order

```

1 Function P-Merge( $L_1, L_2$ ):
2   if  $L_2 = \emptyset$  then  $L \leftarrow L_1$  ;
3   else
4      $m_1 \leftarrow \lfloor L_1.length/2 \rfloor$ ;
5      $L_{1,L} \leftarrow \{L_{1_i}\}_{1 \leq i < m_1}$ ;  $L_{1,R} \leftarrow \{L_{1_i}\}_{m_1 < i \leq L_1.length}$ ;
6      $m_2 \leftarrow$  index where  $m_1$  would be in  $L_2$ ;
7      $L_{2,L} \leftarrow \{L_{2_i}\}_{1 \leq i < m_2}$ ;  $L_{2,R} \leftarrow \{L_{2_i}\}_{m_2 \leq i \leq L_2.length}$ ;
8      $L_L \leftarrow$  spawn P-Merge( $L_{1,L}, L_{2,L}$ );  $L_R \leftarrow$  P-Merge( $L_{1,R}, L_{2,R}$ );
9     sync;
10     $L \leftarrow$  concatenate  $L_L, L_{1_{m_1}}, L_R$ ;
11  end if
12  return  $L$ ;
13 end

```

The complexity is now $\Theta(n)$ while the parallelism grows to $\Theta(n/\log^2 n)$.

Designing efficient algorithms requires:

- A good understanding of the problems
- An advanced knowledge of the underlying mathematics
- Much time spent on a trial and errors approach
- The consideration of all the corner cases

