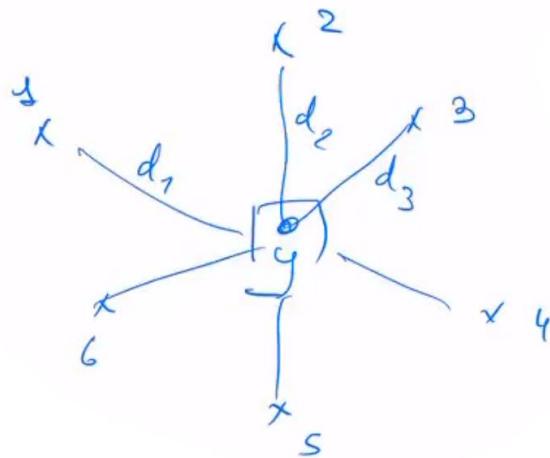


## 1 - Optimization problems and the search space

N persons in a 2D space who want to meet at some location  $y$  which minimizes the total distance travelled by all N persons. What is this point  $y$ ?

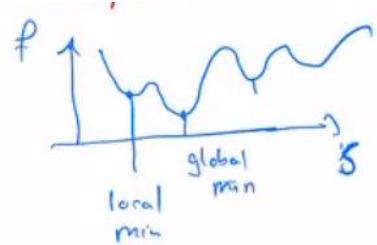


- The search space  $S$ , is a set of possible values for  $y$
  - If the problem is in  $R^2 \rightarrow$  the search Space is  $R^2$ , all different possible values for  $y$
  - The search space can be continuous discrete finite or infinite
  - If the search space is too big, we can't do an exhaustive search
- 
- Lets define a function  $f: S \rightarrow R$ , this function is called objective function, cost function, energy function, fitness function ..
  - Its goal is to quantify a possible value  $y$  of the search space  $S$ , this way we can compare 2 different possible values for  $y$  and choose the best one. The main goal being to find the best value possible for  $y$ .
  - This optimal value will either minimize or maximize the function  $f$

$$\left\{ \begin{array}{l} x_{opt} = \underset{x \in S}{\operatorname{argmin}} / \underset{x \in S}{\operatorname{argmax}} f(x) \\ f_{opt} = \underset{x \in S}{\min} / \underset{x \in S}{\max} f(x) = f(x_{opt}) \end{array} \right.$$

- This  $x_{opt}$  may not be unique

- One difficulty for these problems will be to distinguish global optimums from local optimums, visually it looks like this:



- Usually we deal with multidimensional problems:

$$x \in S \quad x = (x_1, x_2, \dots, x_n) \quad \begin{matrix} \text{component of } x \\ \text{this is a vector} \end{matrix}$$

Here n is called the problem size

This is the number of degree of freedom in the problem.

S

metaheuristics is the way to explore S when no polynomial algorithms exist to find the optimal solution

## 2 - Main principles of metaheuristics, neighbourhood, movements, exploration operator, population metaheuristics

GOAL of metaheuristics: explore a large search space in a clever way. It is needed to solve large problems, for which only exponential algorithms are known

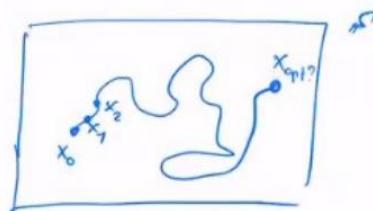
- It is a compromise between CPU time and quality of solution
- No guarantee on the quality of the solution
- As opposed to heuristic where the algorithm is specific to a problem, metaheuristic algorithm can be applied to many different problems

Characterization:

- no hypothesis on the mathematical properties of the fitness function
- Require guiding parameters that defines how the space should be explored. The quality of the solution will depend on how good these parameters are
- need a starting point (usually random)
- need a stopping condition (iteration, no more improvement, time, value ...)
- inspired by natural processes (ant system, beehive ...)
- can be parallelisable
- most of them are stochastic, use random numbers to guide a search

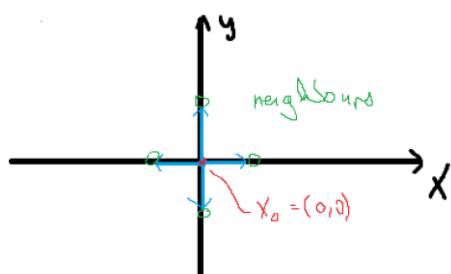
Main principle:

A metaheuristic is an exploration trajectory in the search space  $S$



To move across the search space, we use the concept of neighbourhood.

- If we are currently exploring a point  $x$ , the next point we explore is a neighbour of  $x$
- The neighbourhood is found by doing elementary transformations on the current explored state. neighbourhood of  $x$  is called  $V(x)$



$x_0 \rightarrow x_1 \in V(x_0) \rightarrow x_2 \in V(x_1) \rightarrow \dots$  until the  
 ↑  
 initial solution  
 stopping condition

Search operator:

- When we are at a state  $x$ , we find its neighbourhood  $V(x)$ , and now we need to choose which neighbour to explore next.
- This is done via the search operator  $\overset{u}{\rightarrow} : V(x) \rightarrow y \in S$

$x_0 \overset{u}{\rightarrow} x_1 \in V(x_0) \overset{u}{\rightarrow} x_2 \in V(x_1) \overset{u}{\rightarrow}$   
 ↑  
 initial solution

### Population metaheuristics

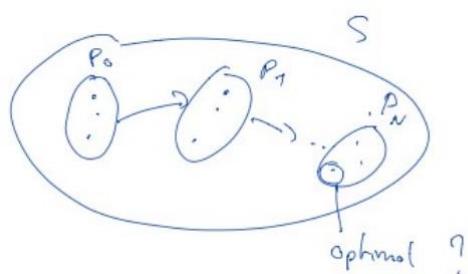
So far we assumed that at each iteration we consider only one possible solution:

$x_0 \rightarrow x_1 \in V(x_0) \rightarrow x_2 \in V(x_1), \dots$

But we could also consider a population of solutions:

$$P_0 = \{x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)}\}$$

M possible  
solution at  
iterat. 0



### 3 – The space of permutations

$S$  is the set of permutation of  $n$  objects

$$n=3$$

$$S = \{(1, 2, 3), (1, 3, 2), (2, 1, 3), (2, 3, 1), (3, 1, 2), (3, 2, 1)\}$$

$$|S| = n! = 6$$

What kind of transformation can we define on a permutation

The simple way is to use transpositions or swap of two entries in the permutation.

$T_i = (i \ j)$  which will swap entry  $i$  with entry  $j$

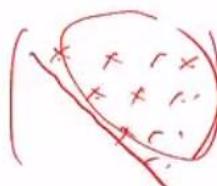
$$\text{For instance: } T_{(2,3)}(1,2,3) = (1,3,2)$$

How many transposition do we have?

$$n=5$$

$$T_i \in \{(1, 2), (1, 3), (1, 4), (1, 5), (2, 3), (2, 4), (2, 5), (3, 4), (3, 5), (4, 5)\}$$

How many are there?  $O(n^2)$



The neighbourhood is of size  $n^2$   
whereas the search space is of size  $n!$

Can we reach any point in  $S$  starting from any initial condition, using a finite number of these transformations?

Answer : Yes see math :

It means that the entire space of permutations is accessible

#### 4 - Complexity and the need for metaheuristics, exploration versus Exploitation

Complexity classes:

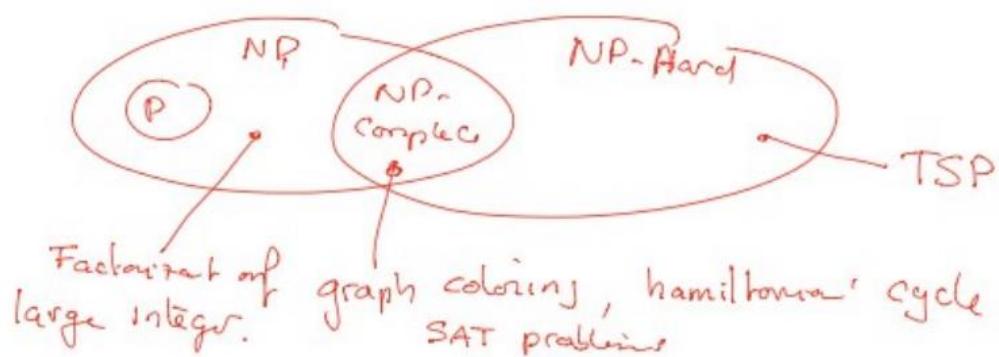
Class P : There is an algorithm that solves the problem in a polynomial time  $T(n) = O(n^m)$

Class NP : Problems for which a solution can be checked in a polynomial time.

Thus  $P \subseteq NP$

Class NP-hard : Those are problems whose solution can be used to solve any NP problem up to a polynomial additional time.

Class NP-Complete Are problems that are both in NP and in NP-hard.



Typically, NP-hard and NP-complete problems can be solved by exponential algorithms

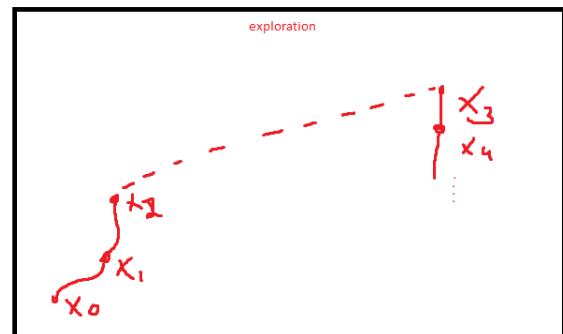
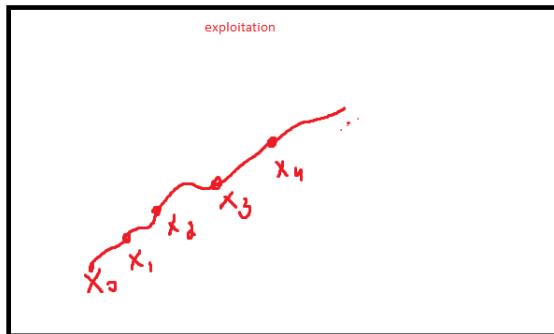
Example:

Hamilton cycle: find a graph that goes through all nodes once and only once can use TSP to be solved.

Exploration vs exploitation:

In all cases, a metaheuristic traverses the search space trying to combine two actions: intensification and diversification, also called exploitation and exploration respectively

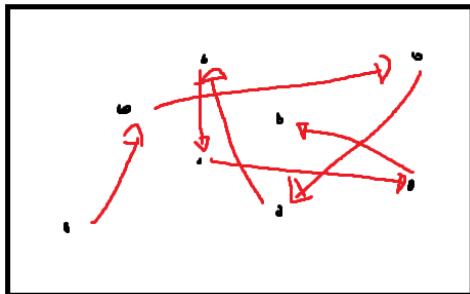
- (exploitation) intensification phase the search explores the neighbourhood of an already promising solution in the search space
- (exploration) diversification a metaheuristic tries to visit regions of the search space not already seen



## 5 - Random search, random walk, hill climbing

Random search:

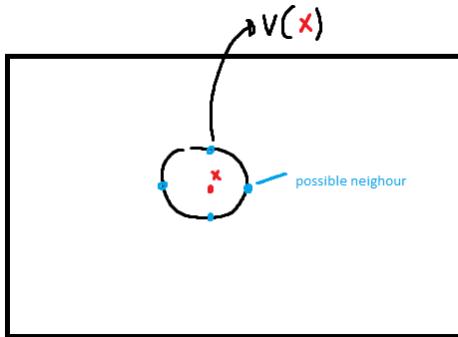
The simplest search method is random search where the next point to test in the search is chosen uniformly at random in the whole search space  $S$ . Usually, one keeps the solution having the best fitness after having performed a prescribed number of steps.



Random walk:

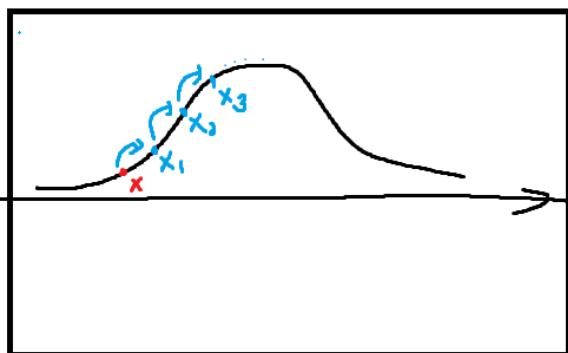
random search to the neighbourhood of the current solution then we have what is called a random walk-in search space.

→ We don't use the operator function to choose which is the best neighbour to explore



Hill climbing:

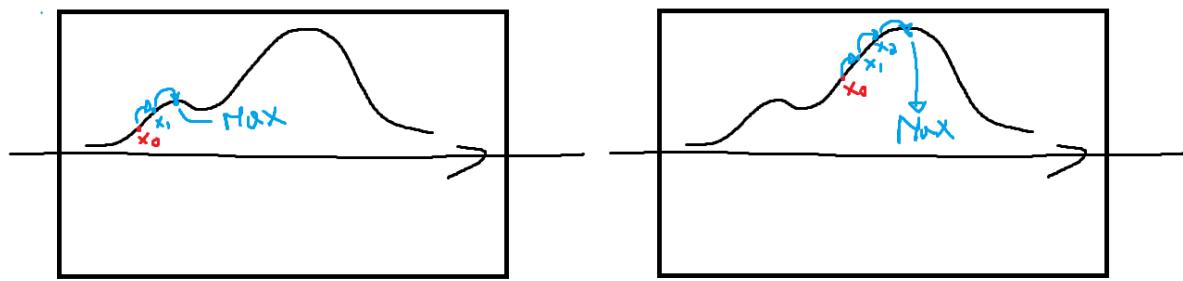
More logical one, we start at a random point  $x$ , look at all our neighbours and choose the one that has the highest (or lowest, depending on the type of problem) fitness to explore next, gives it an effect of climbing (or descending) an hill



We stop when we have climbed the hill -> all neighbours are lower fitness than the current state.

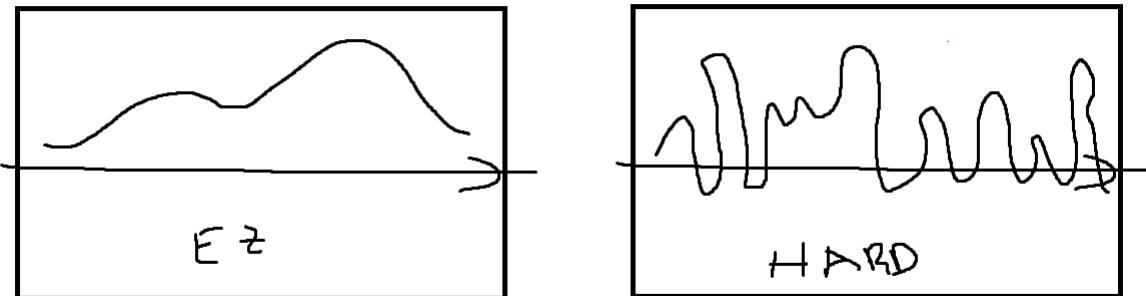
Comment:

Of the three methods the one with the highest chance is the hill climbing. However, it is very easy to get stuck in a local minimum / maximum, it highly depends on the starting point:



as we can see, on the left we got stuck in a local maximum.

Its performance highly depends on the landscape as well



## 6 - NK-problems: motivation, definition, goal

Inspired by biological problems -> genetic regulatory networks

NK problems belong to np-hard class !

Lets give an example of nk landscape problems:

- We consider N persons or agents, labelled with an index  $i, i = 1, \dots, N$
- Each agent  $i$  acts according to two possible strategies, denoted  $x_i = 0$  or  $x_i = 1$
- The success of an agent depends on the strategy it chooses and the type of relation it has (competition or collaboration) with the other persons it interacts with
- If we assume that each agent  $i$  depends on  $K$  other agents, we may define a function  $f_i(x_1, \dots, x_{i+k})$  which gives the profit resulting from the chosen strategy and that of the connected agents

The problem is specified once the  $f_i$  are given, as well as  $K, N$

We want to maximize

$$f = \sum_{i=1}^N f_i \quad \begin{matrix} \text{The profit of} \\ \text{the population.} \end{matrix}$$

Our main interest with NK-problems is to be able to generate synthetic problems of increasing difficulty. (as  $N$  and  $K$  increase, it become more and more difficult to solve.)

### Example MaxOne problem.

Find a chain of  $N$  bits that maximize the number of ones.

Obvious solution:  $x = (x_1, x_2, \dots, x_n) = \underbrace{111 \dots 1}_{N \text{ times.}}$

$$f_i(x_1, \dots) = x_i \quad K=1$$
$$f = \sum_{i=1}^n f_i = \sum_{i=1}^n x_i$$
$$f = N$$

Another example

$$x = \underline{x_1 x_2 \dots x_n} \quad f_i(x_{i-1}, x_i, x_{i+1}) \quad K=3$$

$f_i$  is max for 111 Then it is a trivial solution  
But if  $f_i$  is large for 101 and small for  
Then the chain  $\overbrace{101010101}$

## 7- Tabu Search: main principles and convergence

Used for quadratic assignment problem and TSP

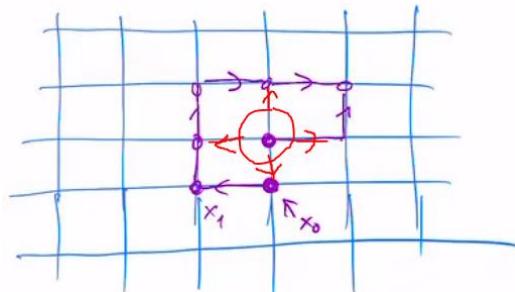
- The search space is explored by going from neighbour-to-neighbour  $x_n \rightarrow x_{n+1} \in V(x_n)$
- The  $x_{n+1}$  is chosen by the search operator, as the non-tabu  $x_{n+1}$  that locally optimizes fitness -> (independent of the current fitness of  $x_n$  -> can be worse) -> if many points have the same fitness, it is random choice
- A tabu list, is a list of already explored states, so that we don't go back
- However, it is possible for a state to leave the tabu list -> after a certain amount of iterations

Tabu list -> prevents us from going back

The fact that  $x_{n+1}$  can have worse fitness than  $x_n$  means we won't get stuck in local max

A state can go out of tabu list, so that the code doesn't get stuck

- This is why we need the list to be updated (if we are in red state, we are stuck):



Convergence:

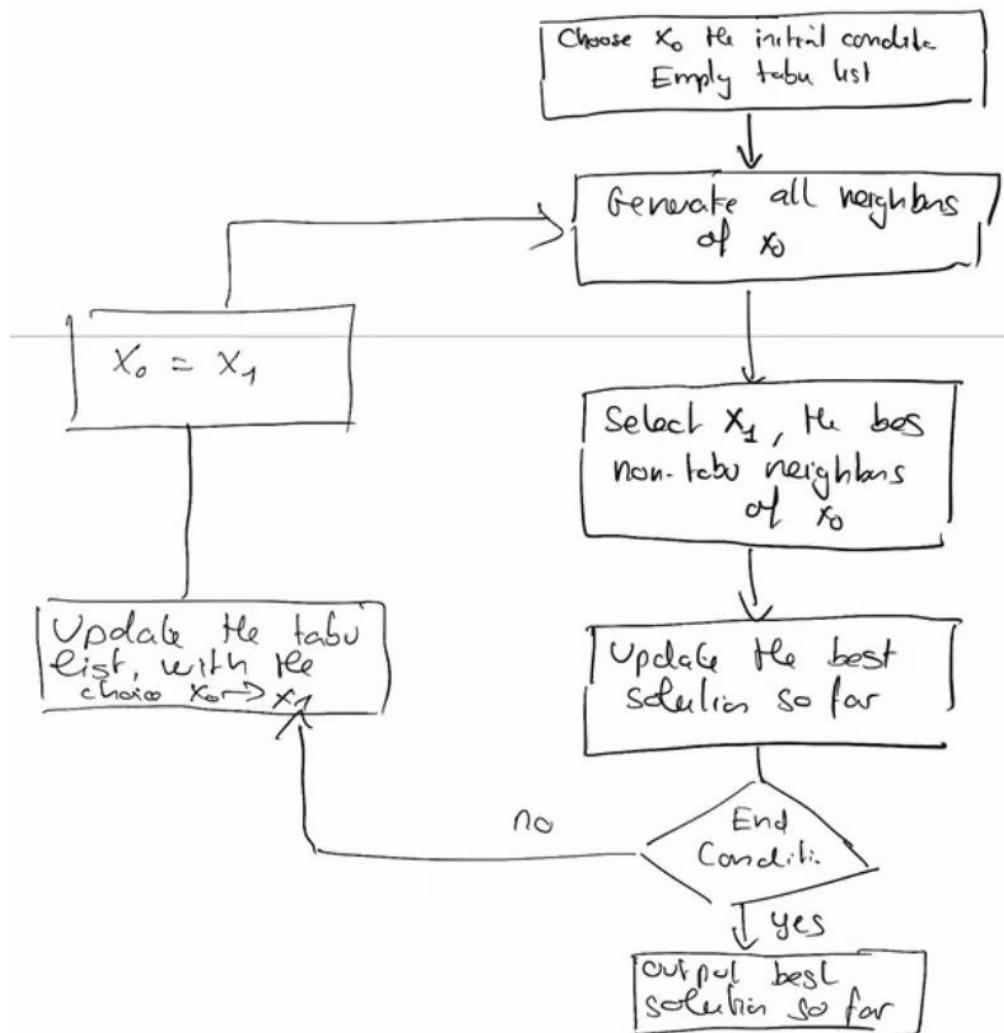
To converge means to find the global optimum!!

-tabu search will converge if:

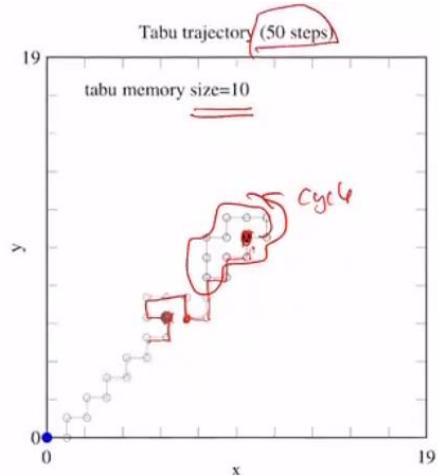
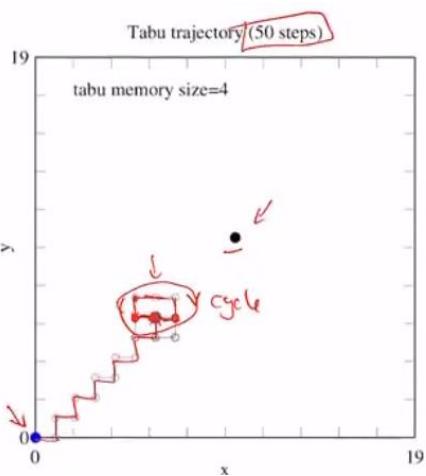
- the search space is finite
- neighbourhood is symmetric:  $x \in V(y) \Leftrightarrow y \in V(x)$
- any state in the search space is reachable by a finite number of steps

- Memorizes all the visited points, but allows the trajectory to use the oldest tabu point to unlock itself -> will visit all the search space -> and find the optimal solution

Flowchart of tabu list algorithm:



Example of the same search space, with 2 different tabu list memories:



## 8- The different ways to implement the tabu list

The goal of the tabu list is to prevent search to explore solution that have been seen previously

There are many versions of a tabu list:

- Visited solutions kept in tabu list for N iterations
- Keep in tabu list all states that have the same fitness as already visited states (if we have visited a state with fitness n, then all neighbours with fitness n go into the tabu list)
- Movements (of transformations) that is forbidden, for example, it is impossible to go left for the next 2 moves -> TO avoid going back to x if we just moved from x to x', the inverse movement is in the tabu list

Updating the tabu list -> adding and removing items: short- and long-term memory

- Short term memory:
    - Use a finite size tabu list, so as we add to the tabu list, the oldest item is deleted
    - Associate a duration to the tabu items (banned time), a movement is banned for N iterations -> N is defined based on knowledge of the problem or by trial and error
  - Long term memory
    - Used to prevent that some movements are never used; the long-term memory can disregard the forbidden movement in the short-term list
    - To implement the long-term memory, we track the movements / states / transitions never done/visited, and at a certain point we force the algorithm to visit them
  - Short list = exploration -> no diversification
  - Long list = exploitation -> with diversification
- We can also put fitness above all, making it ok to revisit a state that has a higher fitness, even if it was in the taboo list

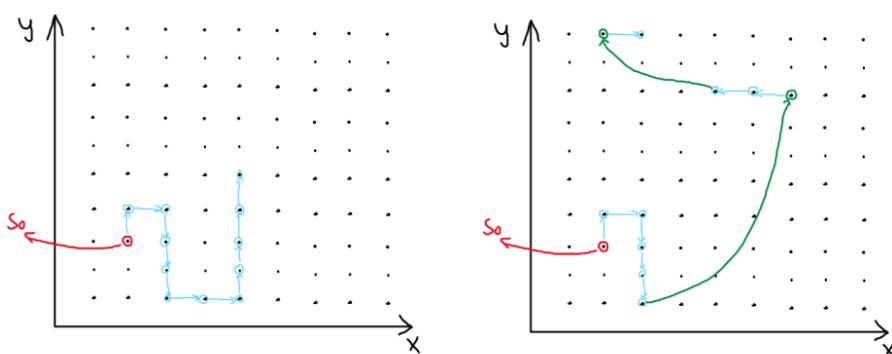


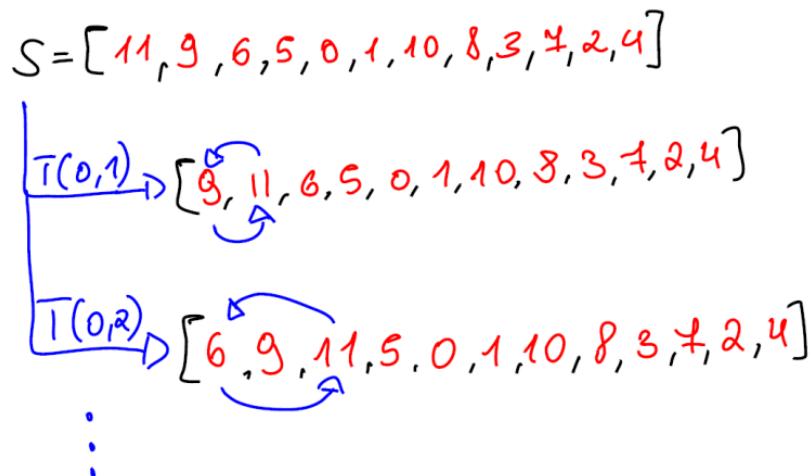
Figure 3: Exploration of  $Z^2$  plane, left without diversification, right with diversification

## 9 - Quadratic Assignment Problems

Combinatorial optimization problem

N objects and n possible locations -> must find the best permutation

- We have n objects and n locations
- We know the distance  $d_{rs}$  between locations r and s
- We are given flows  $f_{ij}$  between all pairs ij of objects
  
- GOAL: find the optimal location for each object -> minimize distance flow -> sum of products

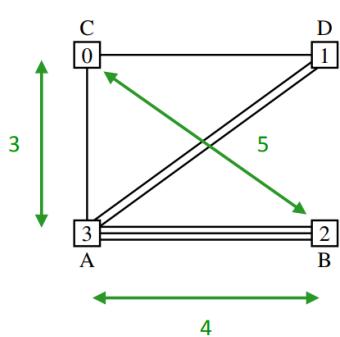


$$f = \sum_{i,j} f_{ij} d_{r_i r_j}$$

- Example: Find the best location (A, B, C, D) for each facility (0, 1, 2, 3) in order to minimize

$$I(\psi) = \sum_{i,j=0}^{n-1} w_{ij} \times d_{\psi_i, \psi_j}$$

distances	flows
$d_{AB} = d_{CD} = 4$	$w_{13} = 2$
$d_{AC} = d_{BD} = 3$	$w_{01} = w_{03} = 1$
$d_{AD} = d_{BC} = 5$	$w_{23} = 3$



Fitness  $I(\psi) = w_{01} \times d_{\psi_0 \psi_1} + w_{03} \times d_{\psi_0 \psi_3} + w_{13} \times d_{\psi_1 \psi_3} + w_{23} \times d_{\psi_2 \psi_3}$

Here  $\psi = (C, D, B, A)$

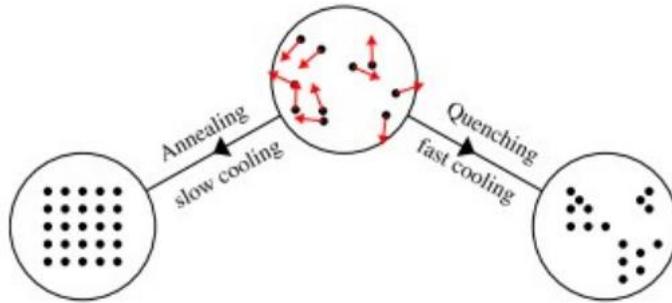
Hence

$$I(\psi) = d_{CD} + d_{AC} + 2d_{AD} + 3d_{AB} = 29$$

## 10 - Simulated annealing, main principles and the Metropolis rule

- Inspired by nature (by physics and metallurgy more specifically)
- Annealing is a process by which a sample is cooled down slowly -> finds global min
- Quench it is a quick cooling down -> finds local min
- Fitness function is called energy function

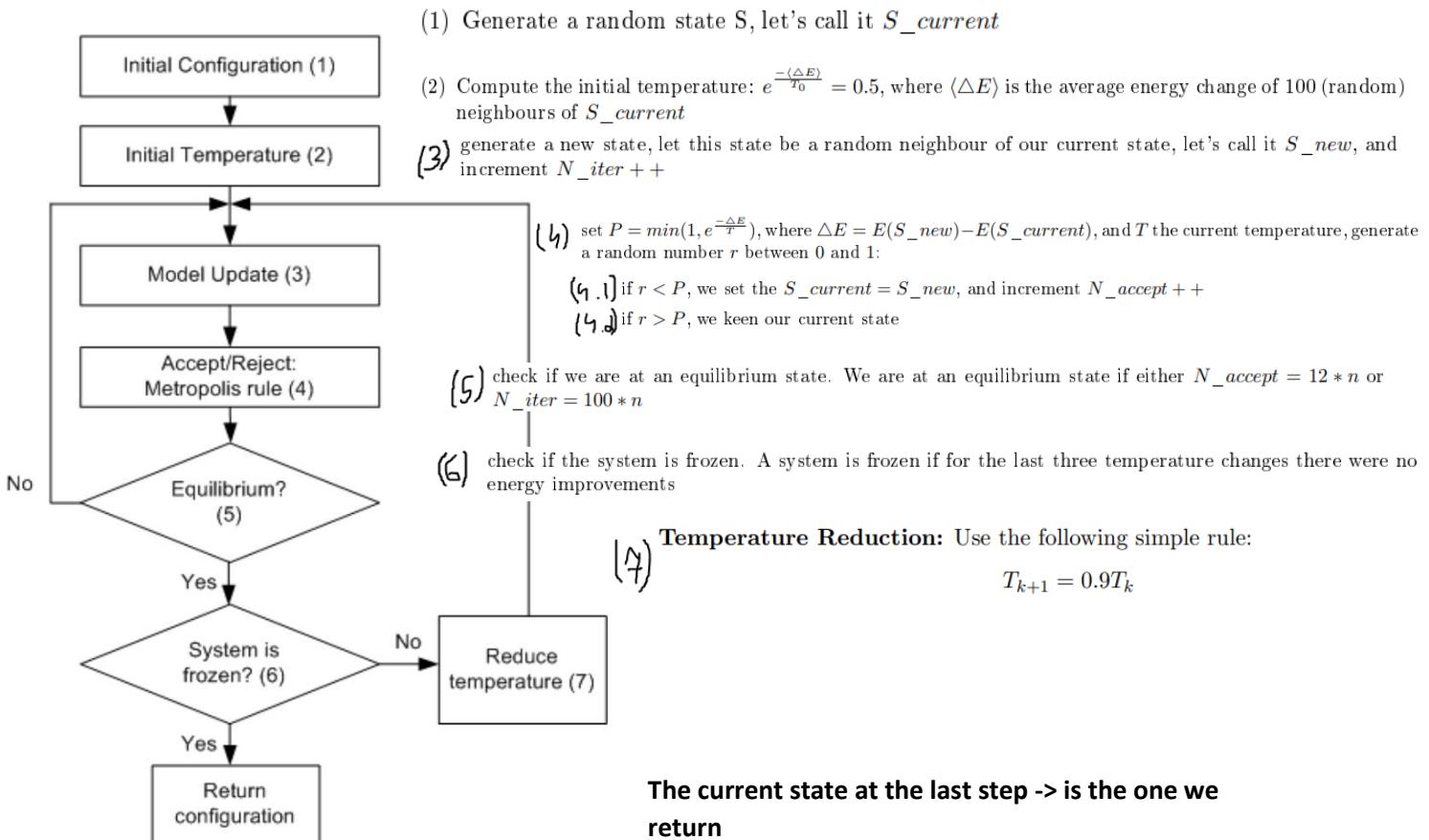
Nature minimizes energy with these processes, and we want to capture this same property with our algorithm



- We start with high temperature, in this state the system explores many possible states
- When the temperature starts to cool down the system is “trapped” does exploitation process

The hope is to have found the global min with the use of both exploration and exploitation

Algorithm:



The current state at the last step -> is the one we return

Metropolis rule:

(4) set  $P = \min(1, e^{-\frac{\Delta E}{T}})$ , where  $\Delta E = E(S_{new}) - E(S_{current})$ , and  $T$  the current temperature, generate a random number  $r$  between 0 and 1:

(4.1) if  $r < P$ , we set the  $S_{current} = S_{new}$ , and increment  $N_{accept}++$

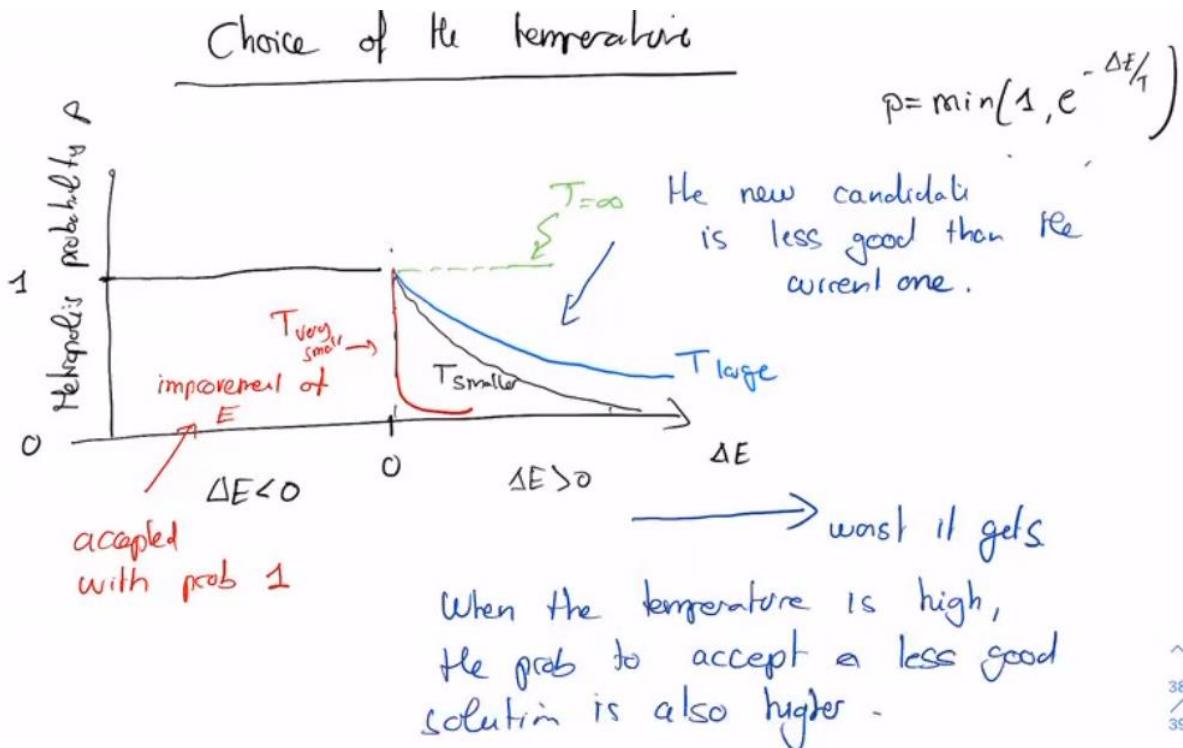
(4.2) if  $r > P$ , we keep our current state

The metropolis rule makes sure we always accept a lower energy state:

- o Because:  $E(S_{new}) - E(S_{curr}) < 0 \Rightarrow E(S_{new}) < E(S_{curr})$
- o And when  $E(S_{new}) - E(S_{curr}) < 0 \Rightarrow P = 1$ , we accept the new state with probability  $P=1$

If the state is of worse energy:

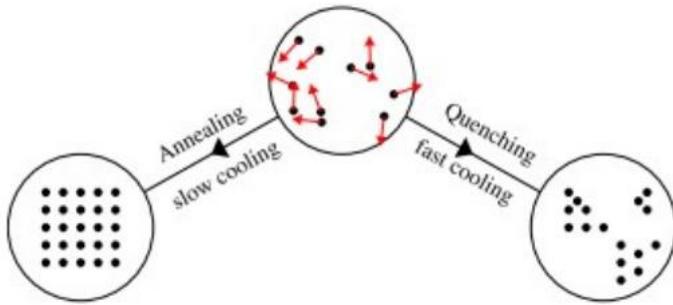
- o Because:  $E(S_{new}) - E(S_{curr}) > 0 \Rightarrow E(S_{new}) > E(S_{curr})$
- o Then there is a chance we accept it still, the higher the temperature, the higher chance to accepting it



## 11 - Simulated annealing, flow chart and choice of parameters

- Inspired by nature (by physics and metallurgy more specifically)
- Annealing is a process by which a sample is cooled down slowly -> finds global min
- Quench it is a quick cooling down -> finds local min
- Fitness function is called energy function

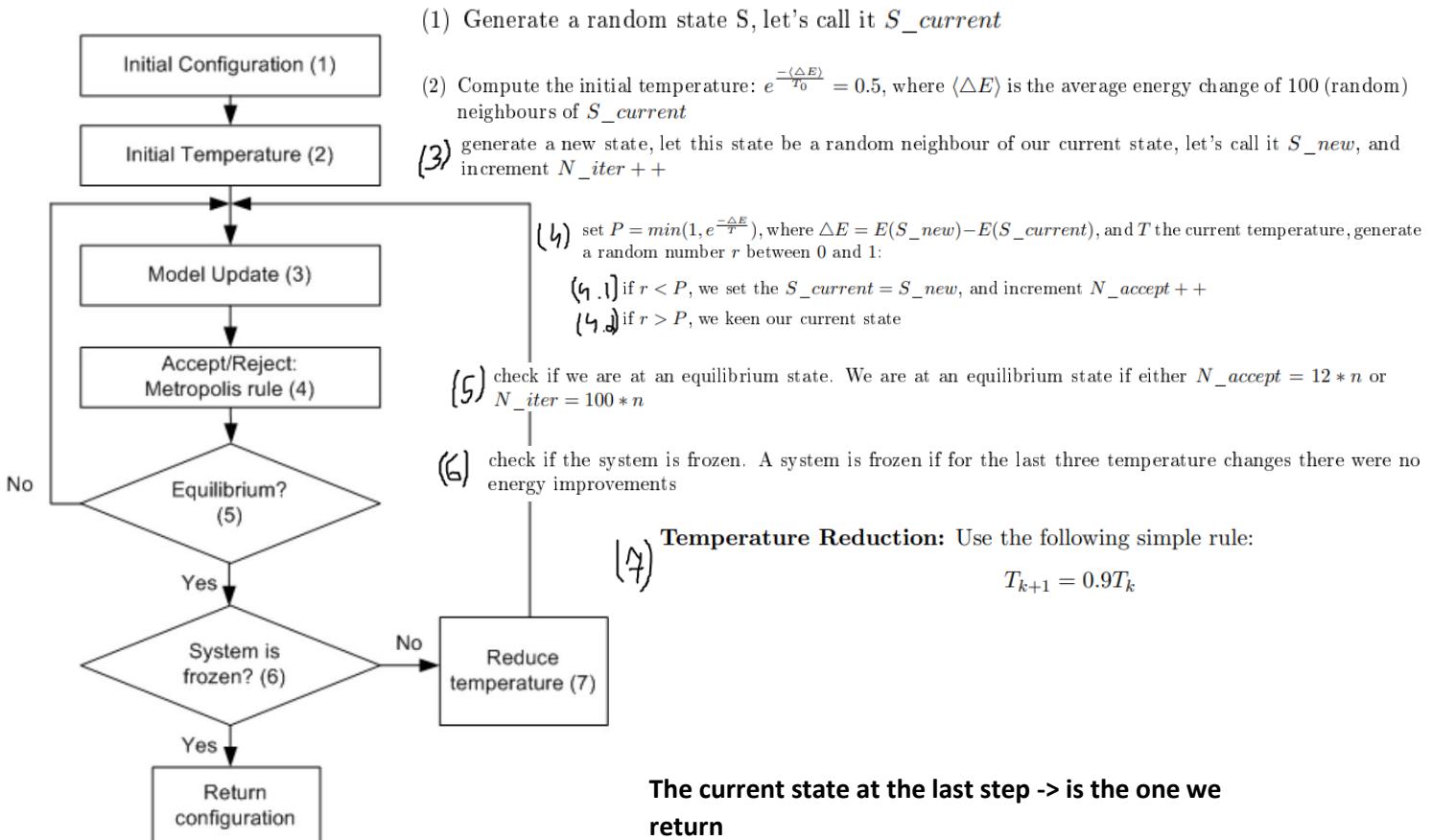
Nature minimizes energy with these processes, and we want to capture this same property with our algorithm



- We start with high temperature, in this state the system explores many possible states
- When the temperature starts to cool down the system is “trapped” does exploitation process

The hope is to have found the global min with the use of both exploration and exploitation

Algorithm:



choice of parameters:

equilibrium:

at step (5), to check if we are at an equilibrium, we check that we have tried  $100*n$  states, or accepted  $12*n$  states -> depending on these values it will be faster or slower have less good or better results

system frozen:

a system is frozen – stop condition – if for the last 3 temperatures there was no fitness improvement -> can increase the number or decrease

temperature:

$tk+1 = tk * 0.9$  -> can lower 0.9 to be in quenching (faster cooling) -> usually worse results

initial temperature might be too low if our random initial state has a low energy

## 12- Convergence of simulated annealing: how to formulate it, how to compute it.

→ Does simulated annealing find the global optimum and under which conditions?

Simulated annealing can be analysed mathematically, it converges in probability.

It gives a solution arbitrary close to the global optimum with a probability arbitrary close to 1

→ The conditions are:

- Movements are irreversible
- Any point in the search space can be reached from any other point in a finite number of movements
- The initial temperature must be large enough (do exploitation before exploration)
- The temperature should not decrease too fast, at iteration t:

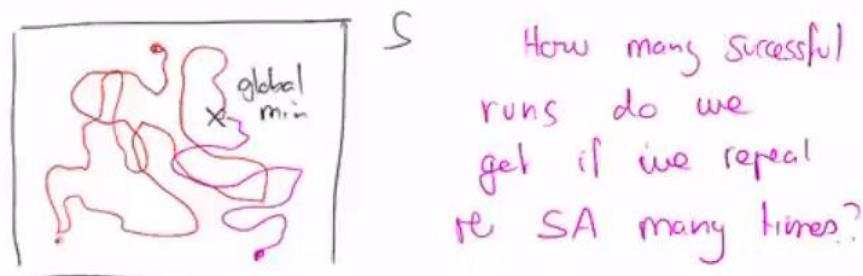
$$T(t) \sim \frac{C}{\log t} \quad \text{for } t \gg 1$$

C is an unknown constant that reflects the variation of E in S.

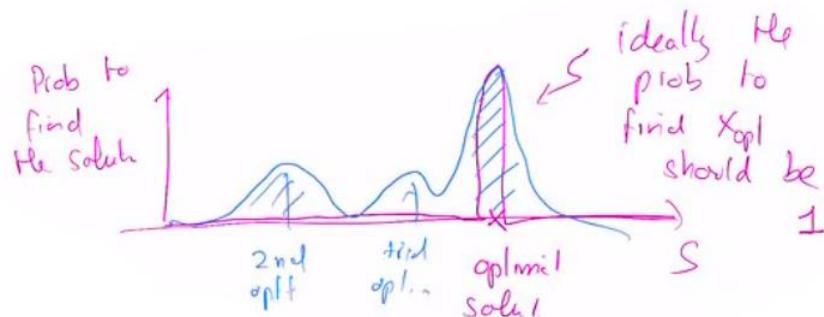
!!! If temperature is decreased too fast, we risk not converging (quenching) !!!

→ In practice,  $\log t$  is too slow. So, temperature will decrease faster

How to compute it:



If we get 100% of success there is convergence, in practice it may end up sub optimal



We want to compute the probability  $P(t, i)$   
that the SA is at point  $i \in S$  at iteration  $t$

$$P(t+1, j) = \sum_i P(t, i) W_{ij}(+)$$



$W_{ij}$  is the transition probability  
from  $i$  to  $j$ . Since we have

100 possible values of  $i$  and  $j$ ,  $W_{ij}$  is a  
100x100 matrix.

$$\begin{matrix} \uparrow & \uparrow \\ |S| & |S| \end{matrix}$$

For SA, we have:

$$W_{ij} = \begin{cases} 0 & \text{if } j \text{ is not a neighbor of } i \\ \frac{1}{k_{\text{out}}(i)} \times P_{\text{metropolis}}(E_i, E_j, T(+)) & \text{otherwise} \end{cases}$$

$\uparrow \qquad \qquad \qquad \uparrow$   
number of neighbors of  $i$        $\min(e^{-(E_j - E_i)/k_B T})$

$k_{\text{out}}(i) = 4$   
in our example.

We also add the following

$$W_{ii} = 1 - \sum_{j \text{ neighbor of } i} W_{ij}$$

## 13- Ant-like algorithms: swarm intelligence, the pheromone trail, and observations about real ants.

These are algorithms that try and mimic effects we see in nature

- Ants as a group can solve difficult tasks, like finding the best path between the ant colony and a food source

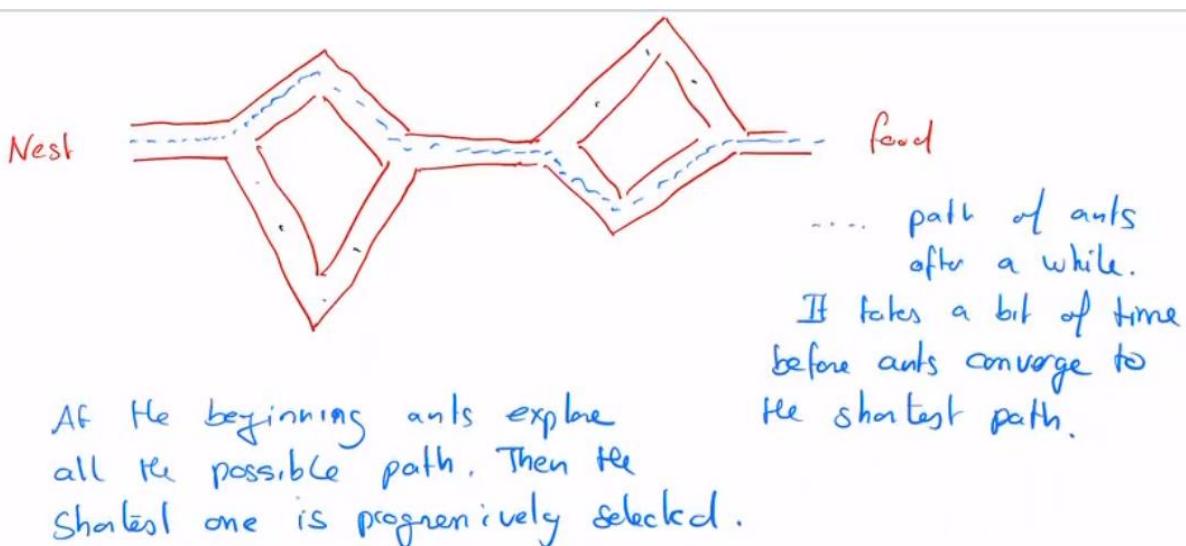
- The important properties of swarm intelligence are:
  - Collaboration is the key to the group's success (the intelligence of the group is more than the sum of the individuals intelligence)
  - This effect of swarm intelligence is also called complex systems, which are a collection of simple entities that interact
  - It results in an emergent behaviour which cannot be understood by any single entity
  - Such behaviours are visible through large scale spatial and temporal organization, sometimes called self-organizations
  - !!! there is no central control, every action is spontaneous !!!

Advantages:

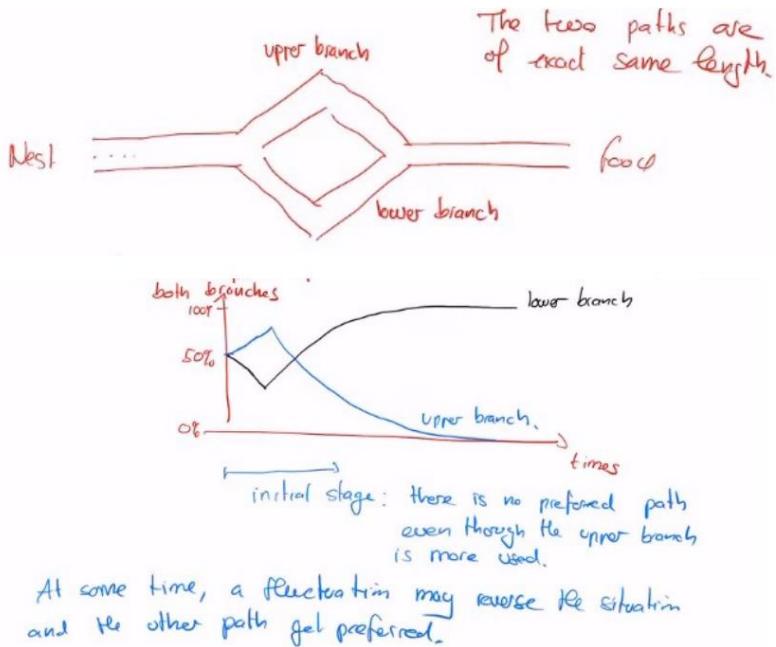
- Robustness to fault tolerance: the process can continue even if an entity disappears
- Possibility of adaptation to new situation
- Natural parallelism in the process, since all entities are 'independent'

Pheromone trail:

- Ants use pheromones to find the shortest path, ants can deposit pheromones to attract other ants to a promising location
- This process is called chemotaxy
- Ants choose the path that smells the strongest
- Pheromones will evaporate over time (how ants forget a suboptimal path)



This experiment was run 10 times, and 9 times out of 10, 100% of the ants were on the shortest path. The 10<sup>th</sup> run, only 80% were on the shortest path.



Explanation:

- (1) We start with no pheromones on the path, so every ant chooses a path at random, leaving pheromones in the trail
- (2) At the second run, it uses the pheromones to choose a path, this path then receives even more pheromones, which attracts even more ants, and so one
- (3) At the end they end up with the shortest path

## 14- Ant System: description of the algorithm for the TSP problem

---

### Algorithm 1

---

```

1: for all  $t = 1, \dots, t_{max}$  do
2:   for all ant  $k = 1, \dots, m$  do
3:     choose a city at random
4:     while there exists a city not visited do
5:       choose a city  $j$  according to (1)
6:     end while
7:     mark a path according to (3)
8:   end for
9:   update all paths according to (2)
10:  Keep the best of solutions obtained at last iteration
11: end for

```

---

(1)

$$p_{ij}^k(t) = \begin{cases} \frac{(\tau_{ij}(t))^\alpha (\eta_{ij})^\beta}{\sum_{l \in J} (\tau_{il}(t))^\alpha (\eta_{il})^\beta} & \text{if } j \in J \\ 0 & \text{otherwise} \end{cases}$$

First off,  $J$  is a set of cities not yet visited in the ant's path, so only cities unvisited can be the next city chosen.

Each unvisited city has a probability of being chosen, this probability is a function of:

(Pheromones in a path to a city \* visibility of the city)/

(Sum of all pheromones\*visibility for all cities unvisited)

(2)

$$\tau_{ij}(t+1) = (1 - \rho)\tau_{ij}(t) + \sum_{k=1}^m \Delta\tau_{ij}^k(t)$$

This function allows us to update the pheromones in each path at each iteration, there is an evaporation rate  $\rho$ , so each path has some of its pheromones that evaporates + an integer that is a function of how many ants took this path -> this value is given by formula (3)

(3)

$$\Delta\tau_{ij}^k(t) = \begin{cases} \frac{Q}{L^k(t)} & \text{if ant } k \text{ used edge } (i, j) \text{ in its tour} \\ 0 & \text{otherwise} \end{cases}$$

$Q$  is a constant,  $L_k$  is the length of the path taken by the ant -> the longer the path -> the smaller is the result -> less pheromones in this path during the next iteration

### Choice of Parameters

- $\alpha = 1, \beta = 5, \rho = 0.1$

- $Q = L_{nn}, \tau_0 = \frac{1}{L_{nn}}$

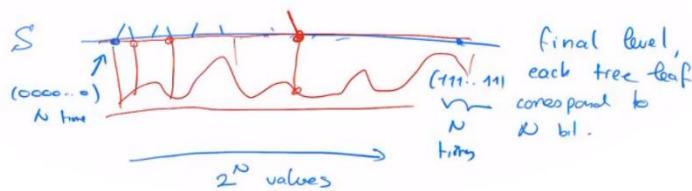
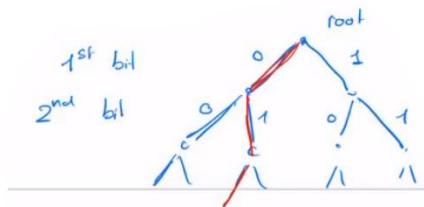
- $m, t_{max}$

## 15- Ant algorithms: the simple version in $\{0,1\}^N$ , and discussion of performance

What we are doing here is a benchmark of this algorithm!

- We will consider a simpler version of the ant algorithm
- We want to find a bit string which maximizes a given fitness function:  $S=\{0,1\}^N$
- This makes it so the search space is a tree structure, at each step (node) we have two options for the next bit, 0 or 1. This tree has length of  $N$

Visually it looks like this:



For every value of the search space  $\rightarrow 2^n$ , there is a fitness function that evaluates the path, then the deposited pheromone on this path is a function of the fitness function result

- At each bifurcation an ant will more likely choose the path that has the most pheromone, each path has the following probability of being chosen:

$$P_{left} = \frac{\tau_{left}}{(\tau_{left} + \tau_{right})} \quad P_{right} = 1 - P_{left}.$$

Statistiques sur 10 répétitions, chacune avec 6 fourmis, après 6 itérations.

effort computationnel	$6 \times 6 = 36$	we explore 36 config. rat=
taux de succès :	<u>7/10</u>	success rate
taux de succès à <u>3%</u> de l'optimum	<u>9/10</u>	

What about a random search?

Let us take 36 points of  $S'$ , at random. What is the probability to find the optimal value within those 36 solutions?

$P = 1/64$  is the prob to find the optimal value at random. With 36 attempt, the prob to find it is

$[0.43] = 1 - (1 - p)^{36}$  is the prob to find it at least once.

## 16 - Particle Swarm Optimization: algorithm and example

- Explores collectively the search space in order to find the optimal solution
- An individual will, at every iteration, choose his path based on his own best path, his current path, and the group's best path
- The hope is that the group of particles will find the best solution

### ALGORITHM

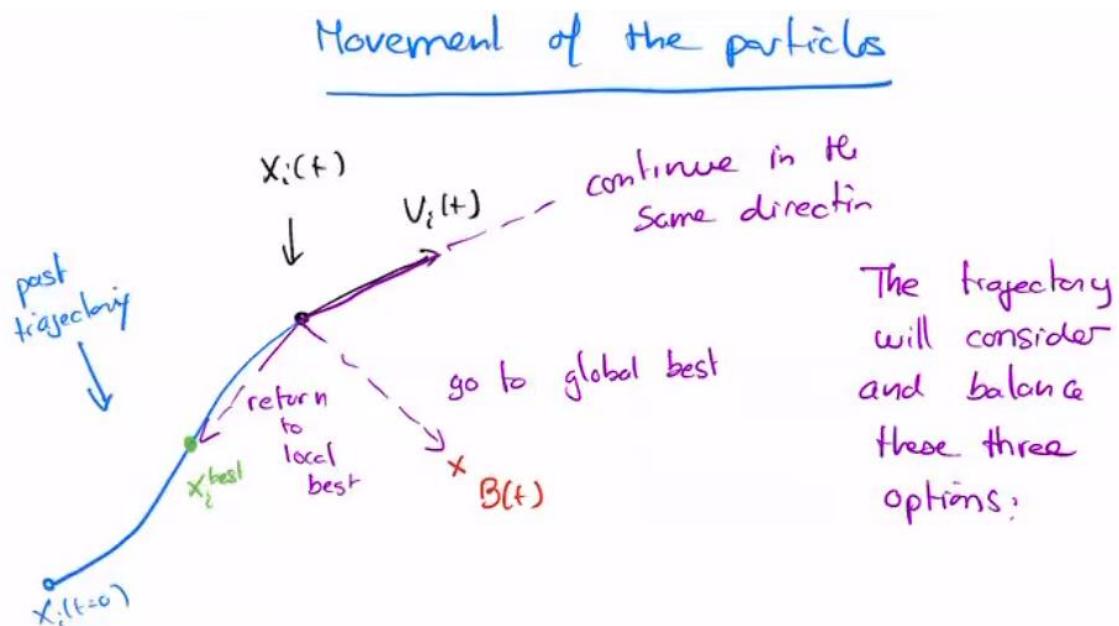
- It is a population metaheuristic (at each iteration the number of individuals stays constant)
- $x_i(t)$  is the position of particle  $i$  at iteration  $t$
- Each  $x_i$  is a possible solution and it has a corresponding fitness
- Particles explore the search space because they have a velocity, which makes them move from iteration  $t$  to iteration  $t+1$

At each iteration each particle will update its local best (his own best solution)

$$x_i^{\text{best}}(t) = \begin{cases} x_i(t) & \text{if } f(x_i(t)) \text{ is better than} \\ & f(x_i^{\text{best}}) \\ x_i^{\text{best}}(t-1) & \end{cases}$$

At each iteration the global best is also updated  $\rightarrow B(t)$

$$B(t) = \underset{x_i^{\text{best}}}{\operatorname{argmax}} f(x_i^{\text{best}}) \quad (\text{for a maximization problem})$$



$$\begin{aligned}
 & \text{current} \quad \text{local best} \\
 \left\{ \begin{array}{l} V_i(t+1) = \boxed{\omega V_i(t)} + C_1 r_1(t+1) \left[ \overset{\text{best}}{x_i}(t) - x_i(t) \right] \\ \quad + C_2 r_2(t+1) \left[ B(t) - x_i(t) \right] \\ x_i(t+1) = x_i(t) + V_i(t+1) \end{array} \right. \quad \text{Global}
 \end{aligned}$$

- $\omega$  is a parameter or  $\omega < 1$  is an inertia parameter.  
How much of the previous velocity we will keep
  - $C_1$  is called the cognitive coefficient as it takes into account the knowledge of the individual.
  - $C_2$  is called the social coefficient as it considers the knowledge of the group.  $C_1 = C_2 \approx 2$
  - $r_1$  and  $r_2$  are random numbers uniformly distributed in  $[0, 1]$
- Everything is a vector!!  
 → Particles are placed randomly in S with  $v = 0$   
 → We need to have search boundaries, many options are available, like a particle bouncing etc...  
 → We maximize V with  $V_{max}$

This algorithm must be applied to a problem that defines velocities addition multiplication etc etc

We can use PSO for computing weights for neural networks

initialize all positions randomly → compute the current fitness for every particle → check if the current fitness is better than the local best for each particle

## 17 - The Firefly algorithm

- Inspired by PSO
- Fireflies attract a mate or prey by emitting light, the higher intensity the higher the attraction is
- For continuous optimization, but exists for discrete

### ALGORITHM

Each firefly  $i$ , at iteration  $t$ , is at location  $x_i(t)$  in the search space:  $x_i(t) \in S \subseteq \mathbb{R}^d$

Each firefly emits a light with an intensity  $I_i(t)$  which depends on the fitness of solution  $x_i(t)$

For a maximization problem, we can simply say

$$I_i = f(x_i(t)) \text{ where } f \text{ is the fitness.}$$

At each iteration, the fireflies move according to:

- One considers all pairs  $(i, j)$  of fireflies with  $1 \leq i \leq n$ ,  $1 \leq j \leq n$ , where  $n$  is the number of fireflies
- If  $I_i < I_j$  then firefly  $i$  moves towards firefly  $j$ , according to the observed attractivity  
to be defined

Attractivity:

$$\exp\left(-\left(\frac{r_{ij}}{\gamma}\right)^2\right) \in [0, 1]$$

where  $r_{ij}$  is the distance separating fireflies  $i$  and  $j$  and  $\gamma$  is a parameter weighing this distance.

This quantity will be the fraction of the distance  $r_{ij}$  that the less intense firefly will move towards the more intense one.

Depending on how we define  $\gamma$  this attractivity will affect far away fireflies, depending on how it is tuned it will be exploration vs exploitation

Movement:

Let us assume that  $i$  moves towards  $j$  each component  
then

$$x_i^{\text{new}} = x_i + e^{-(f_i - f_j)^2} (x_j - x_i) + \alpha \left( \frac{\text{rand}() - \frac{1}{2}}{d} \right)$$

$\uparrow$   
fine for  
continuous optimization.

parameter  
 $\alpha \in [0, 1]$

$\in [-\frac{1}{2}, \frac{1}{2}]$

this a d-dimensional  
vector of random number  
between 0 and 1

Code:

- Initialize the fireflies position randomly in the search space
- for every couple  $(i,j)$  of fireflies
- if intensity of  $i <$  intensity of  $j$ , we update position and intensity of firefly  $i$
  
- we don't change firefly  $j$ !!!

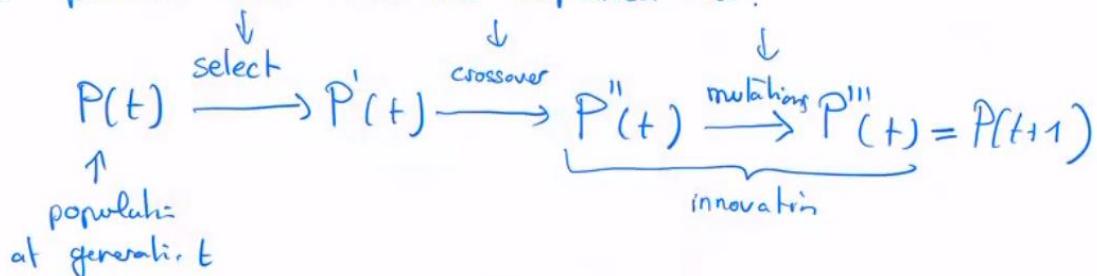
## 18- Genetic algorithms: inspiration and algorithm

- Inspired by Darwin evolution
  - On représente les individus dans S par des chromosomes ou son patrimoine génétique
  - Le degré d'adaptation est donné par la fitness
  - The population evolves at each generation
  - We choose the best individuals over the less fit ones
  - The population size stays constant over generations
- 
- We will mostly consider maximization problems

### ALGORITHM:

- We start with a randomly generated generation
  - The next generation is derived from the previous one by following the steps bellow
- (1) Select the best individuals in the population (many strategies can be applied here)
  - (2) Apply recombination/ crossover among the individuals as well as mutation
- 
- If the genetic process loses the best individual, it is readded to the population by replacing the current worse! Which means we track the global best

The process can also be expressed as:



The selection process focusses on exploitation, and mutation focus on exploration

Crossover does a bit of both -> exploration a bit more

## 19- The different selection operators

There are several ways to implement the selection step

### (1) Fitness proportional

Draws with replacement, we choose randomly n times an individual from the current population. The individual is chosen proportional to its fitness

More chance to choose a high fitness individual

There can obviously be copies of individuals

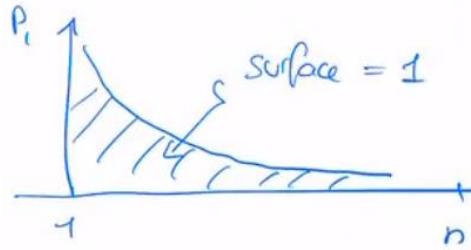
Il faut faire attention à s'il y a des individus avec une fitness négative

$$P_i = \frac{f(s_i)}{\sum_{k=1}^n f(s_k)}$$

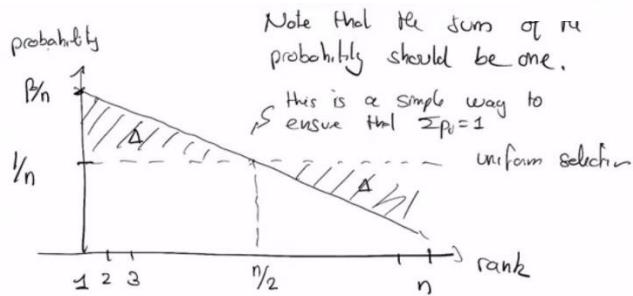
### (2) Selection by rank

We set the individuals in an order depending on their fitness value (by rank), this works with negative fitness results, as well as for minimization problems

For each rank we assign a probability of the individual being chosen, by having higher probability for higher ranked individuals



### (3) Selection by linear rank



Thus is called linear rank selection and the function above is

$$P_i = \frac{1}{n} \left[ \beta - 2(\beta-1) \frac{i-1}{n-1} \right] \quad i=1, \dots, n.$$

$$\begin{cases} \beta > 1 \rightarrow P_i > P_n \\ \beta < 2 \rightarrow P_n > 0 \end{cases}$$

### (4) Selection by tournament

k-tournament selection -> draw k random elements from the population, each individual having  $1/n$  probability of being chosen

from the k chosen individuals, select the one with the highest/lowest fitness depending on the problem

do this process n times

→ One doesn't need to define a fitness function, just needs to be able to compare individuals

## 20- The takeover time in genetic algorithms: goal, definition, value for tournament selection.

Quantify how fast a selection mechanism eliminates all but the best individuals

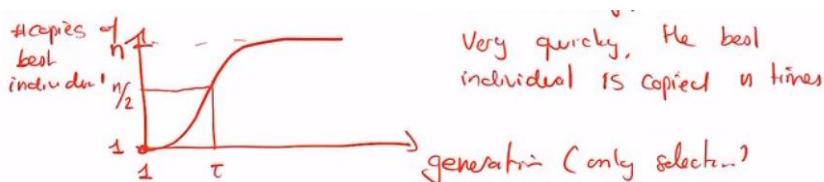
Let's define  $\tau$  the takeover time as the number of selection rounds (no crossover and mutation), to have only the best individuals

It measures the intensity of selection:

- $\tau$  Is small -> strong selection
- $\tau$  Is large -> weak selection

"large" and "small" is proportional to size of population  $n$  -> we will show that

$$\tau = O(\log n)$$



→ Selection is often exponential

Theory:

→ The number of copies  $m(t)$  of the best individual at generation  $t$  goes as:

$$\frac{dm}{dt} = \alpha m \left(1 - \frac{m}{n}\right)$$

↓  
prop. to the current number

↑  
limitate due to a finite population of size  $n$

This solution for tournament selection:

At each selection we have a probability of choosing the best individual:

$(1 - \frac{m}{n})^k$  ← repeat k times  
 ↑ prob to take the best  
 ↓ prob not to take the best

As well as the probability of choosing the individual at least once:

$$1 - [1 - (1 - \frac{m}{n})^k]$$

$$m_{t+1} = n \left(1 - \left(1 - \frac{m}{n}\right)^k\right)$$

$$\geq n \left(1 - \left(1 - \frac{m}{n}\right)^2\right)$$

$$\text{with } t = \ln(n-1) \Rightarrow m(t) = \frac{n}{2}$$

$$= n \left(1 - \left(1 - \frac{2m}{n} + \frac{m^2}{n^2}\right)\right)$$

$$= n \left(\frac{2m}{n} - \frac{m^2}{n^2}\right)$$

$$\approx m \left(2 - \frac{m}{n}\right)$$

To have half of the population filled with the best individual, one has to perform  $O(\log n)$  selections.

## **21- Genetic Programming: goal, fitness evaluation, function and terminal Sets**

The goal is to apply evolutionary ideas to computer programs

- Make the code evolve to solve the problem
- We will consider a population of “computer programs”, that will be selected, recombined and subject to mutation
- Until we reach the desired solution
- What can be expected of such programs? -> they should discover a function that produces the right output for a given input

FITNESS:

*output, input*

Let us consider a program  $p$  such that:  $y = P(x) \rightarrow x$  and  $y$  can be complex data structures

Let  $A$  be a training set:  $A = \{<x_1, y_1>, <x_2, y_2>, \dots, <x_k, y_k>\}$

For input  $x_i$  we know the output  $y_i$

Which means the fitness will be computed by:  $f(P) = \sum_{i=1}^k |y_i - P(x_i)|$

FUNCTION SET AND TERMINAL SET:

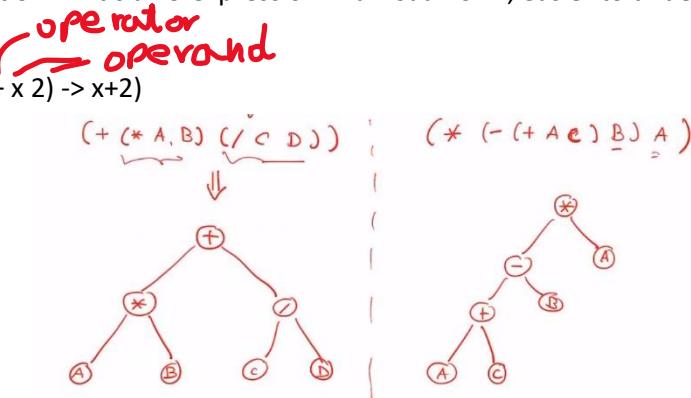
The space of programs that we can build is specified by 2 sets:

- $F$ : set of nodes functions  $\{+, *, -, /\}, \{\text{AND}, \text{OR}, \text{NOT}\}$  -> these are the operators that modify the operands
- $T$ : set of terminal nodes  $\{A, B, C, D\}$  -> these set contains the variables and constants involved in the problem

## 22- Genetic Programming: tree representation, initialization, crossover, mutation, and bloat.

Tree representation: -> it's an S-expression in a visual form, easier to understand the order of operations

(S-expression:  $(+ x 2) \rightarrow x+2$ )



### INITIALIZATION:

- As always, the population is generated randomly with values from the terminal as well as function set. With this tree representation, the terminal sets are the leaves, and the function set is the nodes
- We can set a probability  $p$ , to draw from F set, and  $1-p$  to draw from T set
- Probability  $p$  will influence the depth of the tree
- It is wise to have a dynamic probability as we go deeper in the tree, at the beginning favouring a draw from F set, and later favouring a draw from T set

### CROSSOVER:

- We exchange at random the sub-trees of the parents at random

### MUTATION:

- We chose a random node and rebuild randomly the sub-tree from this node
- We can change the mutation probability depending how deep we are in the tree

### BLOAT:

By doing crossover and mutation our functions may become more and more complex

In practice we should prevent this effect by having max length parameters in place

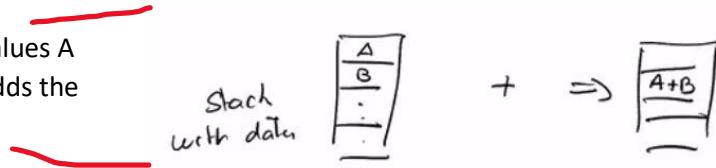
Sometimes it is possible that a subtree can be simplified such as  $(- x x)$

## 23- Genetic Programming: the stack-based, instruction-driven representation.

The idea is to get closer to a procedural programming language -> set of instructions executed sequentially

To make sure that a program stays valid after the genetic transformation we use a **stack based** approach

The **+** operator consumes both top values A and B, performs the operation, and adds the result to the top of the stack



If we code as such, we will be able to have individuals that all have the same length

This will be called **instruction driven representation**

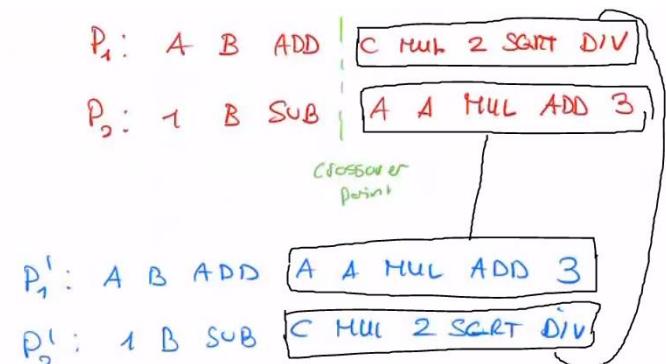
$$\text{ex: } A \ B \text{ ADD } C \text{ MUL } 2 \text{ SQRT } \text{ DIV} = \frac{(A+B)C}{\sqrt{2}}$$

However, it is possible that at some point there aren't enough values in the stack to do an operation (1), or that we end with more than 1 value on the stack (2), if this happens, we:

- (1) Skip this operation and leave the stack as is, do the next instruction
- (2) Return the value on top of the stack

MUTATION AND CROSSOVER:

We cross the strings of instructions with each other for crossover, and mutation is done with a mutation probability



- on peut ajouter des branchements IF-ENDIF

$I_1 \dots I_k \text{ IF } I_{k+2} \dots I_n$        $\text{END IF } I_{n+2}$   
   $\xrightarrow{\geq 0? \text{ oui}} \quad \xrightarrow{\text{non}} \rightarrow \text{flag}\{0,1\}$

- on peut ajouter des boucles

$I_1 \dots I_k \text{ LOOP } I_{k+2} \dots I_n \text{ ENDLoop } I_{n+2}$

## 24- Evolution Strategy: individual and population versions.

Single individual metaheuristic  $x \in \mathbb{R}^d$  evolving by mutation only

Population version also including crossover

Metaheuristic for continuous optimization  $S \subseteq \mathbb{R}^d$

### (1+1) ALGORITHM -> 1 parent and 1 child

- The child is produced by a gaussian mutation of the parent:
- $X'$  is the child,  $x(t)$  the parent, and  $N(0, \sigma)$  is the mutation
- The mutation follows a normal distribution
- The child is only accepted if its fitness is better than the parents
- The  $\sigma \in \mathbb{R}^d$  which means  $\sigma$  is also a vector
- $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_d)$
- Each dimension can have a different  $\sigma$
- (1+1)-ES is much like a random walk search with hill climbing strategy
- If we are accepting too many children, we increase  $\sigma$ , to make larger jumps
- If we accept too few children, we decrease  $\sigma$
- We should want to be around 1/5 acceptance rate

$$x' = x(t) + N(0, \sigma)$$

$$X(t+1) = \begin{cases} x' & \rightarrow f(x') > f(x) \\ x & \rightarrow f(x) < f(x') \end{cases}$$

### POPULATION ALGORITHM

The size of the population (nb of solutions) is denoted  $\mu$

Each individual is represented as  $(x^i, \sigma^i) \quad i=1, \dots, \mu$

$x^i \in S, \sigma^i$  is the associated mutation parameter

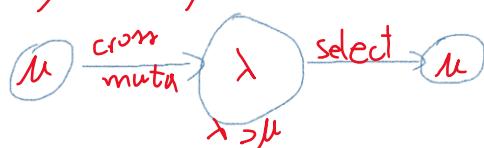
There are 2 variants of the population approach:

①  $(\mu + \lambda)$ -ES:  $\lambda$  children generated from  $\mu$  parents



②  $(\mu, \lambda)$ -ES: in this case  $\lambda > \mu$  children are generated from the  $\mu$  parents,  $\mu$  are selected

This last selection is deterministic, we choose the  $\mu$  best individuals



### CHILDREN GENERATION

#### MUTATION

- Choose 2 parents among the  $\mu$  possible parents
- Apply mutation to that child
- Repeat  $\lambda$  times to generate  $\lambda$  children

$$\left\{ \begin{array}{l} x^1 = x + N(0, \sigma) \\ \sigma^1 = \sigma + e^{N(0, \Delta \sigma)} \\ \Delta \sigma = 1/\sqrt{\lambda} \end{array} \right.$$

#### Crossover

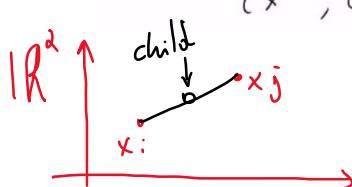
The crossover is performed both on  $x$  and  $\sigma$

• The uniform crossover is the following

$$(x_e^{\text{child}}, \sigma_e^{\text{child}}) = (x_e^{\text{parent}_1 \text{ prob } \frac{1}{2}}, x_e^{\text{parent}_2 \text{ prob } \frac{1}{2}}, \sigma_e^{\text{parent}_1 \text{ prob } \frac{1}{2}}, \sigma_e^{\text{parent}_2 \text{ prob } \frac{1}{2}})$$

or one can also use the arithmetic crossover. Let us consider two parents,  $i$  and  $j$

$$(x_e^{\text{child}}, \sigma_e^{\text{child}}) = \frac{1}{2}(x^i, \sigma^i) + \frac{1}{2}(x^j, \sigma^j)$$

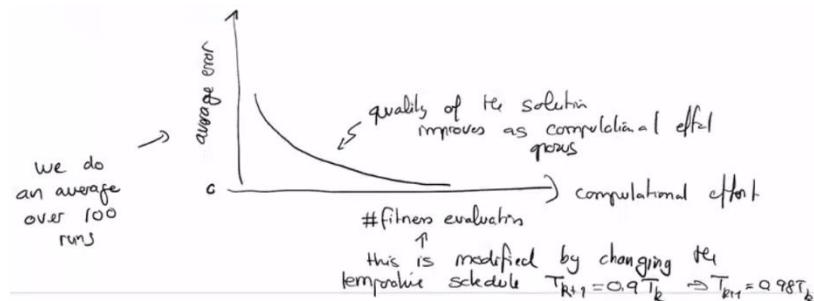
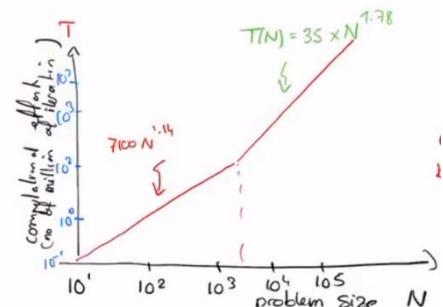


## 25- Performance of metaheuristics: examples, specificities of the performance evaluation, metrics, approach, "No Free Lunch" theorem

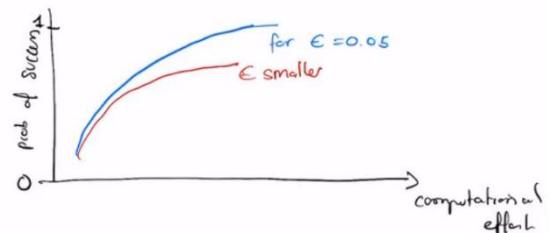
- No guarantee of quality of solution
- They are stochastic, so the behaviour is always different
- To determine the success of an algorithm in a problem, we must use statistics

EXAMPLE:

- Let's consider SA for TSP, we have  $n$  towns randomly placed in a  $2 \times 2$  spatial region
- We will generate many such problems, varying  $N$
- Our question is: how long does it take does it take for SA to give an answer? (This does not mean that the answer is the global optimum, but when SA stops)
- Performance speed:
  - With  $N \in [20, 2000]$ , the difficulty is linear  $O(N^{1/4})$
  - With  $N \in [5000, 50000]$ , difficulty is  $O(N^{1/7.8})$  which is smaller than quadratic  $O(N^2)$
  - It is also much better than exhaustive search  $O(N!)$
- Performance quality (1):
  - To compute the quality of a solution one must know the best solution for a given problem
  - We will place 50 cities in a circle and see how SA finds an accurate solution



- Performance quality (2):
  - We can also compute the probability of success, out of 100 runs, how many gave us the global optimum? (accuracy of  $\epsilon$ )



The metrics used to evaluate a metaheuristic are:

- The complexity in time to get a solution
- Average error as a function of computational effort
- Probability of success
- Statistics are needed, between 100 and 1000 runs

Question: if metaheuristic A is better than B on a given problem, can we conclude that A is always better?

Formulation of NFL:

Let us consider a finite search space  $S$  of size  $|S|$

We consider fitness function

$$f: S \rightarrow Y \quad \text{where } Y \text{ is a finite subset of } \mathbb{R}$$

All possible problems are then specified by a given  $f$  sampled from a number  $|Y|^{|S|}$  of possibility.

One considers a computational effort  $m$ , meaning that we consider a trajectory of exploration of  $m$  points

$$(x_i, f(x_i)) \quad i=1, \dots, m$$

- A metaheuristic cannot be better than another one on all possible problems.
- For any performance metric, no algorithm will be better if all discrete fitness function are considered.
- If  $A$  is better than  $B$  on a given class of problems, there is another class where  $B$  will be better.

-  $m$  iterations

### NFL Theorem

Let  $P(d_m | f, m, A)$  the probability that trajectory  $d_m = \{x_1, f(x_1), \dots, x_m, f(x_m)\}$  generated by metaheuristic  $A$  contains the optimal value of  $f$ .

$$\sum_f P(d_m | f, m, A) = \sum_f P(d_m | f, m, B)$$

Thus, on average all metaheuristics behaves the same when compared on all possible problems.

But in practice, not all possible problems have the same probability and some are highly pathological and not realistic.

## 26- Phase transition in optimization problems: problem description and properties.

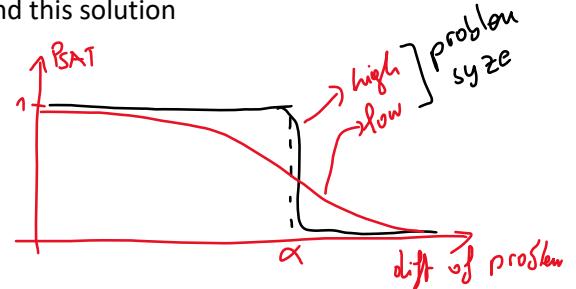
We will analyse statistically a metaheuristic for solving problems of increasing difficulty

We will consider a sub-class of satisfaction problems

We will analyse this problem analytically, we will define the probability that a random instance of the problem has a solution and whether the metaheuristic will find this solution

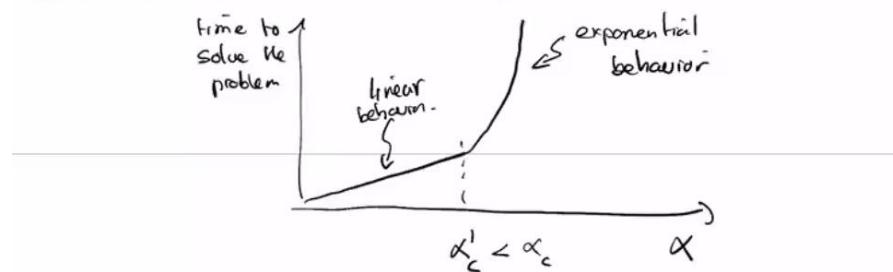
The specificity of the problems we will consider is that:

- if the difficulty parameter is low, the probability that a random instance has a solution is one.
- If this parameter increases up to the critical value the probability drops abruptly to zero.



This abrupt behaviour is called phase transition

We will also see that the metaheuristics used to find solution has another phase transition in the CPU time needed to solve the problem:



**SAT Problem** has N Boolean variables and M Boolean equations

We want to find an assignment of these N variables that satisfy all M equations

The goal is to find an assignment of these N variables that satisfy all M equations, if it exists, we say the problem is satisfiable

These problems can be turned into an optimization problem, the goal is to minimize the energy E, which is the number of unsatisfied equations

→ Example with only XOR operation

$$\begin{cases} x_1 + x_2 + x_3 = 1 \\ x_2 + x_4 = 0 \\ x_1 + x_4 = 1 \end{cases} \quad \begin{array}{l} N=4 \quad (4 \text{ variables}) \\ \quad x_1, x_2, x_3 \text{ and} \\ H=3 \quad x_4 \\ \quad (3 \text{ equations}) \end{array}$$

$$(x_1 \ x_2 \ x_3 \ x_4) = \begin{cases} (1 \ 0 \ 0 \ 0) \\ (\cancel{0} \ 1 \ 1 \ 0) \end{cases}$$

→ This problem has a solution, so the minimal energy is 0, if a problem can't be satisfied, then its energy optimal is 1

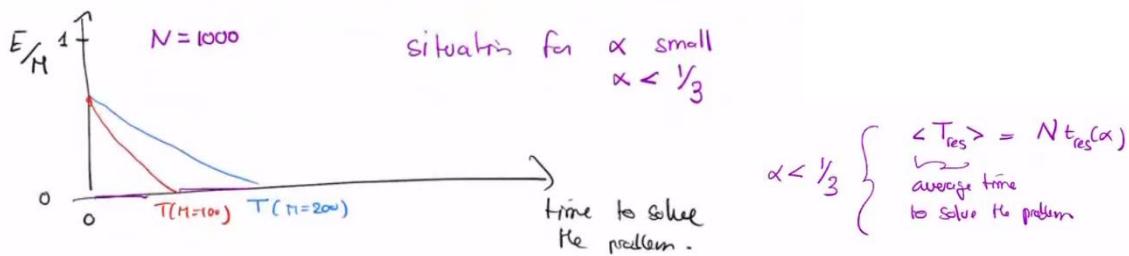
## 27- Phase transition in optimization problems: the RWSAT algorithm and its behaviour.

RWSAT (random walk SAT) can be used to find solution to any SAT problem, in particular XORSAT problems.

Pseudocode:

- N variables, M equations, k variables per equation
- Initialize all N variables at random
- Compute E # of UNSAT equations
- t=0 # number of iterations
- while E>0 and t<tmax:
  - o choose at random one of the non-sat equations
  - o choose at random one of its k variables
  - o set this variable to its complement # 1 -> 0, 0 -> 1
  - o compute E # might have worse energy than before the change
  - o t = t + 1
- print E, t # t is time to solution

How does this RWSAT algo behaves?



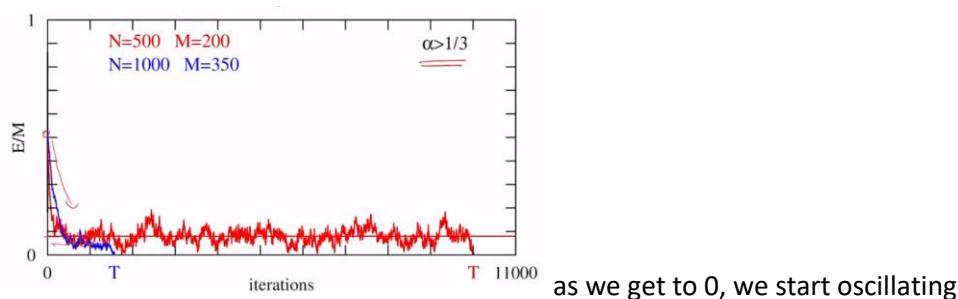
Resolution time is linear with  $N$ , with a coefficient that increases with  $M$

We see that a solution is found ( $E=0$ )  
in a time which increases with  $M$  ( $\alpha = \frac{M}{N}$ )

HOWEVER, when  $\alpha = M/N$  increases, the average execution time becomes exponential

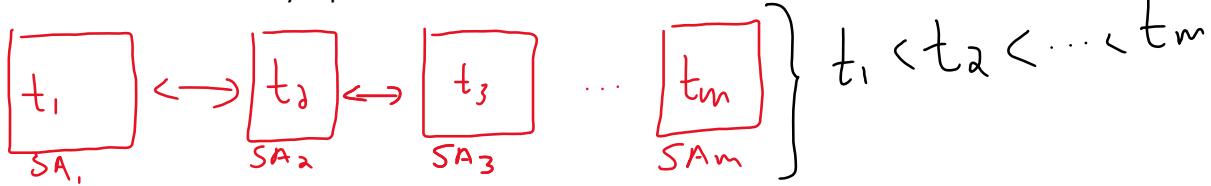
$$\left. \begin{array}{l} \alpha > \frac{1}{3} \\ \text{and of course} \end{array} \right\} \begin{array}{l} \langle T_{\text{res}} \rangle = \exp(N T_{\text{res}}) \\ T_{\text{res}} = T_{\text{res}}(\alpha) \end{array}$$

$$\alpha < \alpha_c = 0.4174$$



## 28- The parallel tempering.

- Parallel version of SA
- We have many replicas of SA



- As opposed to standard SA, here the temperature of each  $SA_i$  is constant (might change just a bit), all these  $SA_i$  run in parallel. They will interact by exchanging configurations (current solution)
- Neighbouring systems can exchange their current solution according to a probability law
- Let's consider 2 systems,  $i$  and  $j=i+1$  (neighbours) with temperatures  $T_i$  and  $T_j$ , and current solution  $C_i$  and  $C_j$ , with energy  $E_i$  and  $E_j$

$$p_{ij} = \min(1, e^{-\Delta_{ij}}) \rightsquigarrow \Delta_{ij} = (E_i - E_j) \cdot \left(\frac{1}{T_j} - \frac{1}{T_i}\right)$$

$$p_{ij} = P_{ji}$$

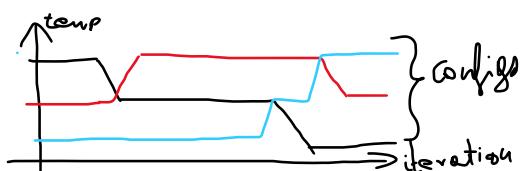
- If  $T_i < T_j$  and  $E_i > E_j \rightarrow \Delta_{ij}$  is negative, and  $p_{ij} = 1$
- What this means, is that if we have a good solution at high temperature, it will always exchange with a less good solution at lower temp
- The opposite is also possible, a good solution in a high temp can also go to a higher temp, with  $p < 1$

! in parallel tempering the temperature schedule is distributed across the replicas of SA, and not over iterations !

- At low temperature we exploit the solution
- At high temp we explore other solutions and other regions of  $S$
- Exchanging configuration is a way to combine diversification and intensification

GUIDING PARAMETERS

It obviously depends on how we define our parameters



- How many SA?  $M$  is often:  $M=\sqrt{N}$   $\rightarrow N$  is the problem size
- At what frequency should we consider exchange of configuration? When both systems reach an equilibrium state, according to the SA definition
- What is the range of temperature between  $T_1$  and  $T_M$ ?  $T_1$  should be small enough to allow the convergence,  $T_M$  should be large enough to allow exploration ( $T_M = T_0$  of standard exploration)
- We can fine tune the temperature of each system over time, if the exchange rate is too large  $> 2\%$  we increase the temperature difference between all systems
- If the exchange rate is too low  $< 0.5\%$  all temperature difference between the systems is decreased