

Optimization for Machine Learning

Lecture 6: Discrete Optimization

December 10, 2020

TC2 - Optimisation

Université Paris-Saclay



Anne Auger and Dimo Brockhoff

Inria Saclay – Ile-de-France

Course Overview

Date		Topic
Thu, 5.11.2020	DB	Introduction to (Continuous) Optimization
Thu, 12.11.2020	AA	Continuous Optimization I: differentiability, gradients, convexity, optimality conditions
Thu, 19.11.2020	AA	Continuous Optimization II: constrained optimization, Lagrangian relaxation, gradient-based algorithms, stochastic gradient
Thu, 26.11.2020	AA	Continuous Optimization III: stochastic algorithms, derivative-free optimization
Thu, 3.12.2020	AA	Discrete Optimization I: graph theory, greedy algorithms Continuous Optimization IV
Thu, 10.12.2020	DB	Discrete Optimization
Thu, 17.12.2020		Final exam

Concrete Information About Exam

Written for those who can be there

- multiple choice, typically 4 answers each (1-4 answers correct)
- closed book (nothing allowed but pen) → easier questions 😊
- next Thursday (Dec. 17) @ ~~1:30pm~~ **1:45pm**
- 2 hours

Oral exam for those who can't be there for the written exam

- also closed book 😊
- 20 min slots via Zoom or MS Teams
- please let me know today if you are one of those students
 - best by e-mail during the break (include your name and your availability)
- we will schedule the exams by tomorrow
- possible slots Thursday or Friday morning next week (optimally all consecutive)

Integer Programming

- variables are integers
- simplest example:
optimization in $\{0, 1\}^n$

ML example:

hyperparameter tuning with
algorithm parts being present
($x_i = 1$) or not ($x_i = 0$)

Combinatorial Optimization

- Search space not necessarily
anymore a subset of \mathbb{R}^n
- for example, optimization on
graphs

ML example:

structure optimization of neural
networks

Exercise: Differences Continuous/Discrete Opt.

What are the differences between continuous and discrete optimization?

optimality conditions

gradient direction?

local/global optima

convexity

neighborhoods

Discrete vs. Continuous Optimization

Important Differences/Observations

- finite search space → still: enumeration impracticable
- discrete neighborhood, sometimes not even clear how to define
- gradient inexistent → follow locally best neighbor?
- different neighborhoods, different definition of local optimum!
example later
- partial evaluations common for discrete problems
- blackbox vs. greybox vs. whitebox
...meaning that solvers for discrete problems are typically more specialized

Overview Discrete Optimization

Algorithms for discrete problems:

- often highly problem-specific
- but some general concepts are repeatedly used:
 - greedy algorithms
 - branch and bound
 - dynamic programming
 - randomized search heuristics

Motivation for this Last Part of the Lecture:

- get an idea of the most common algorithm design principles
- we cannot
 - go into details and present many examples of algorithms
...but for a few
 - analyze algorithms theoretically with respect to their runtime

Greedy Algorithms

Greedy Algorithms

From Wikipedia:

“A *greedy algorithm* is an algorithm that follows the problem solving *heuristic* of making the locally optimal choice at each stage with the hope of finding a global optimum.”

- Note: typically greedy algorithms do not find the global optimum

Lecture Outline Greedy Algorithms

What we will see:

- ① Example 1: Money Change problem
- ② Example 2: ϵ -Greedy Algorithm for Multi-Armed Bandits

Example 1: Money Change

Change-making problem

- Given n coins of distinct values $w_1=1, w_2, \dots, w_n$ and a total change W (where w_1, \dots, w_n , and W are integers).
- Minimize the total amount of coins $\sum x_i$ such that $\sum w_i x_i = W$ and where x_i is the number of times, coin i is given back as change.

Greedy Algorithm

Unless total change not reached:

add the largest coin which is not larger than the remaining amount to the change

Note: only optimal for standard coin sets, not for arbitrary ones!

Related Problem:

finishing darts (from 501 to 0 with 9 darts)

Example 2: Multi-Armed Bandits

- generic problem of resource allocation
- classic reinforcement learning problem showing the exploration–exploitation tradeoff dilemma

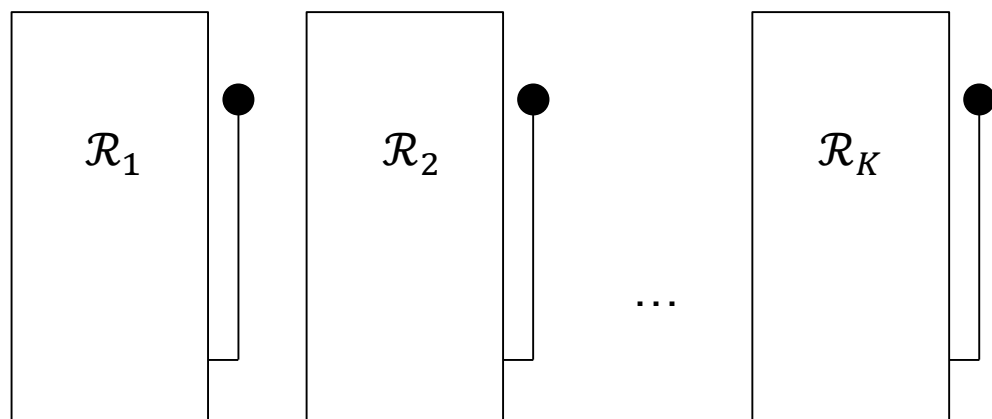


Yamaguchi
i先生

Example 2: Multi-Armed Bandits



Yamaguchi
i先生



- K single-arm bandits with a lever
- Each bandit has a fixed but unknown probability distribution \mathcal{R}_i attached to it with a mean μ_i
- At each time step t , we decide to pull a lever (i) and get a reward r^t according to \mathcal{R}_i
- Overall, we want to maximize the sum of the rewards
- The regret after T steps is defined as $\rho = T\mu_{max} - \sum_{t=1}^T r_t$

Exploration vs. Exploitation: The ϵ -Greedy Algorithm

Exploration: pull new levers (or underexplored ones) to get better estimates on the expected rewards

Exploitation: pull the arm, we think is the best arm

...the latter being the greedy approach here

The ϵ -Greedy Algorithm

- With probability $1-\epsilon$: pull the lever, we think is best
- With probability ϵ : pull a random lever (uniformly)

To be decided (not discussed further here):

How to estimate the probabilities (e.g. pulling each lever once at first)

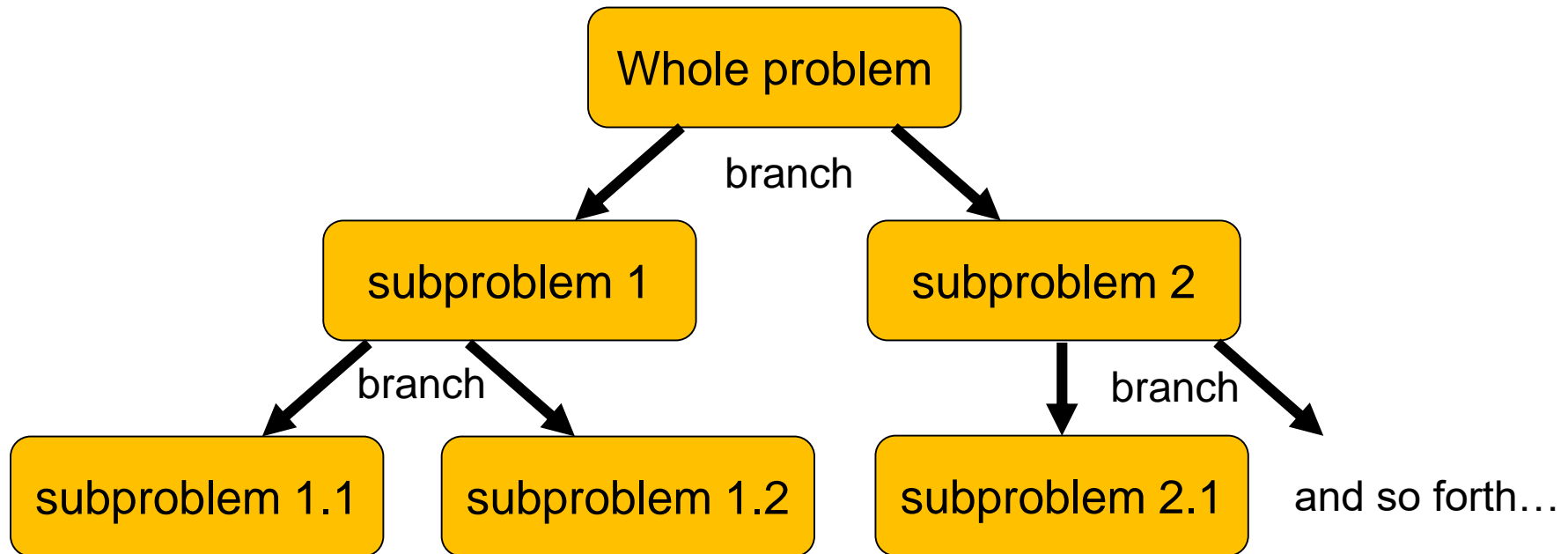
How to choose ϵ (constant vs. decreasing over time)

constant ϵ gives linear regret

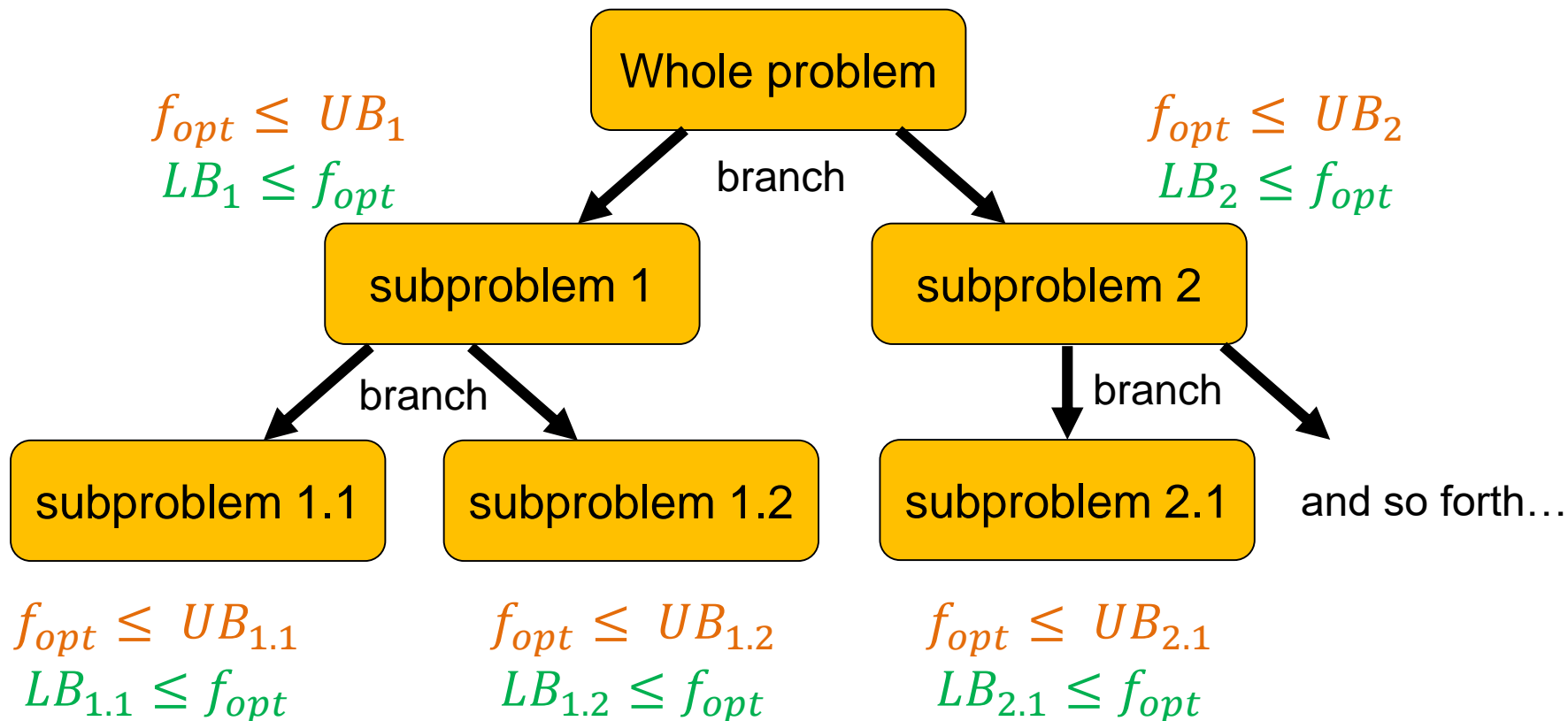
Branch and Bound

Idea Behind Branch and Bound

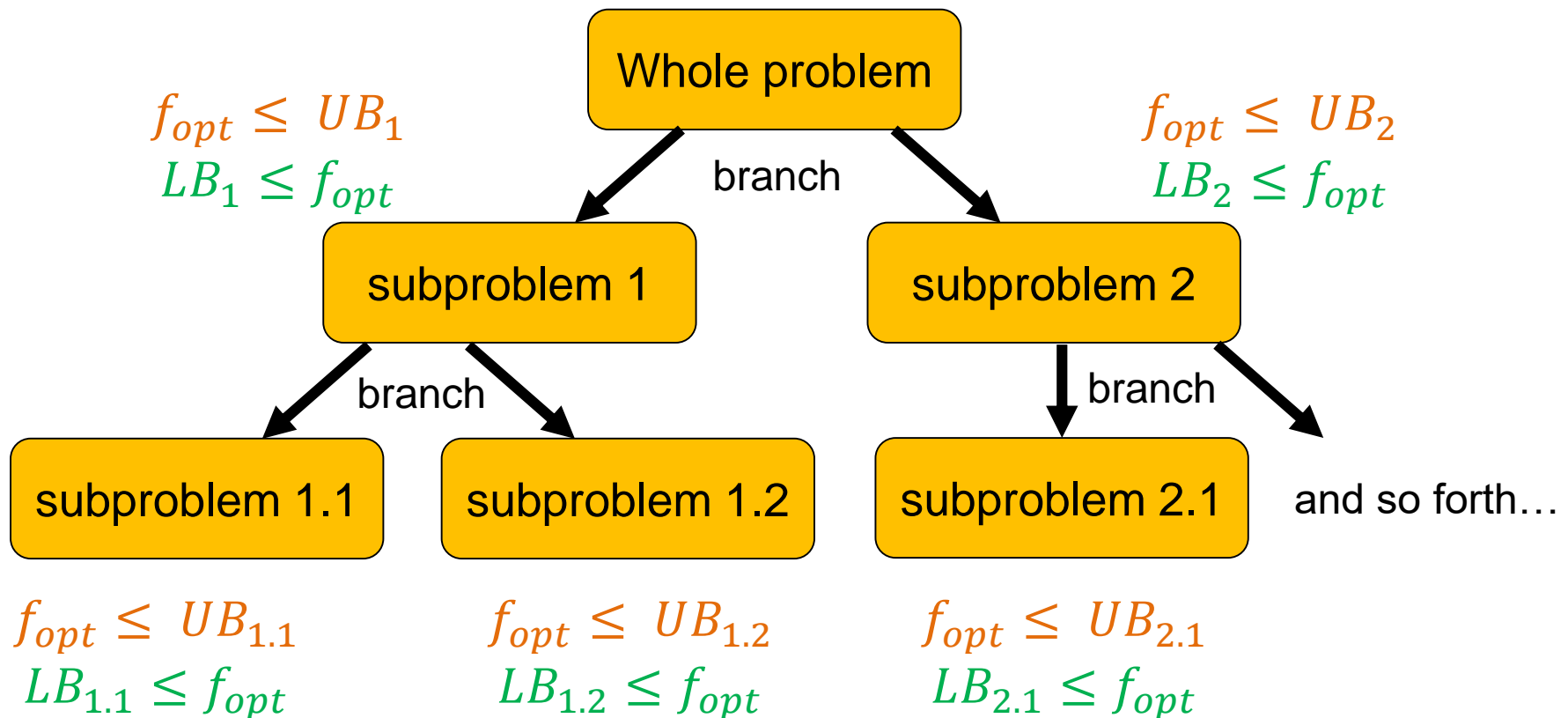
- Basically enumerates the entire search space
- But uses clever strategies to avoid enumerations in bad areas



Idea Behind Branch and Bound

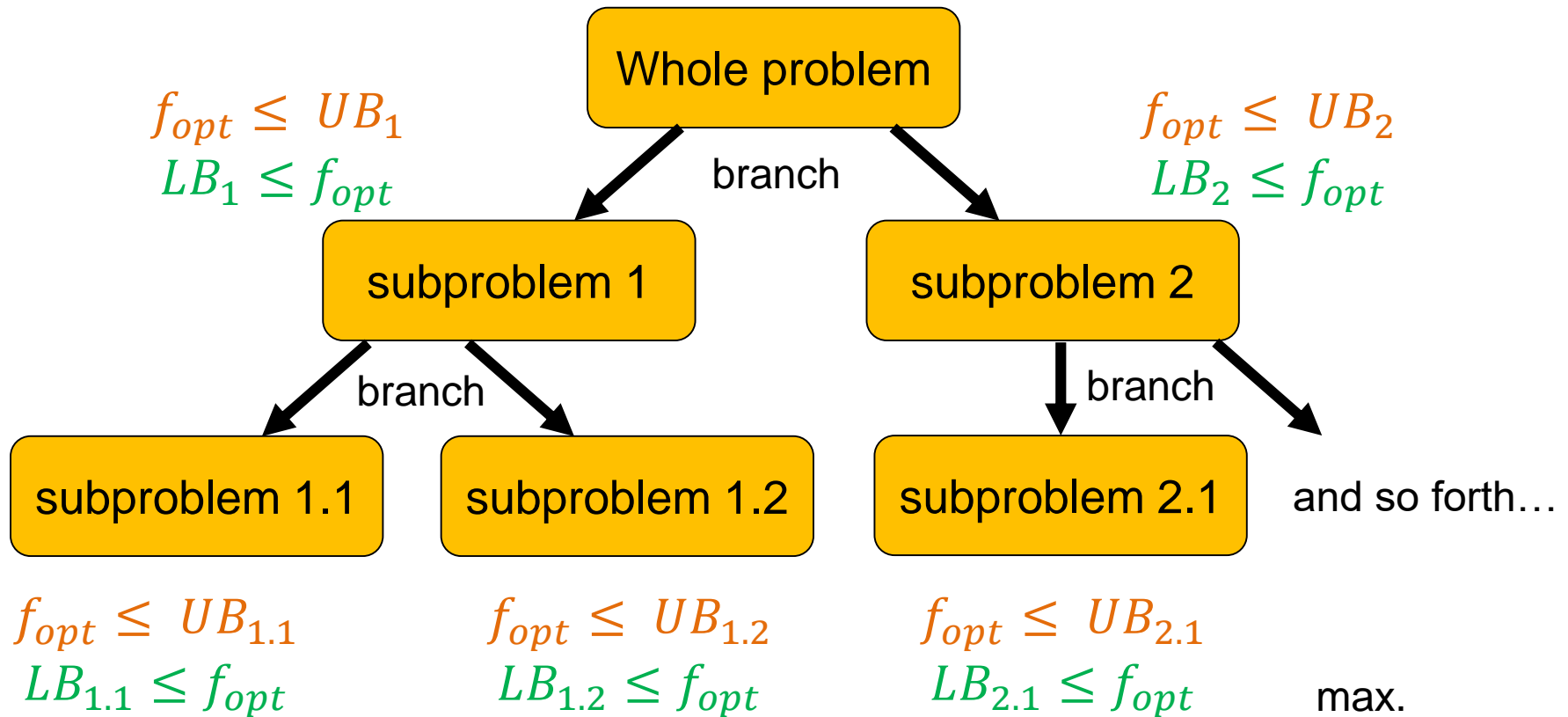


Idea Behind Branch and Bound



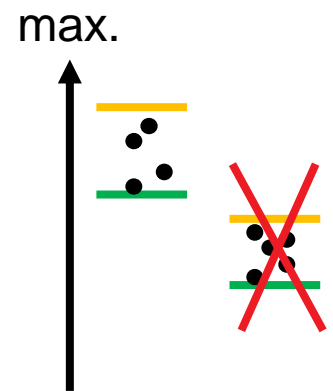
when can we actually avoid evaluating all solutions?

Idea Behind Branch and Bound



We can stop exploring/branching if

- $UB=LB$
- UB for new subproblem lower than LB for another
[when maximizing]



How do we get Upper and Lower Bounds?

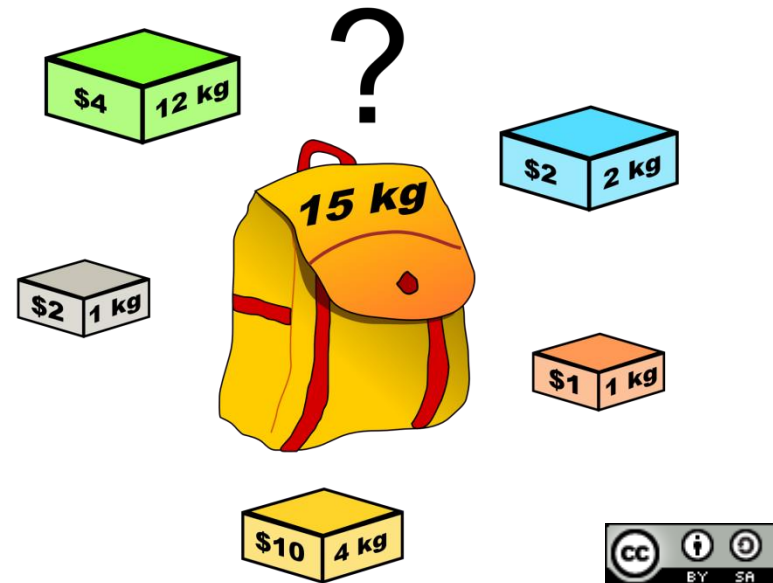
We assume again maximization here...

- A feasible solution gives us a lower bound
the optimum will be at least as good as a solution, we know
- Hence, fast (non-exact) algorithms such as greedy can give us lower bounds
- For upper bounds, we can relax the problem
for example, by removing constraints

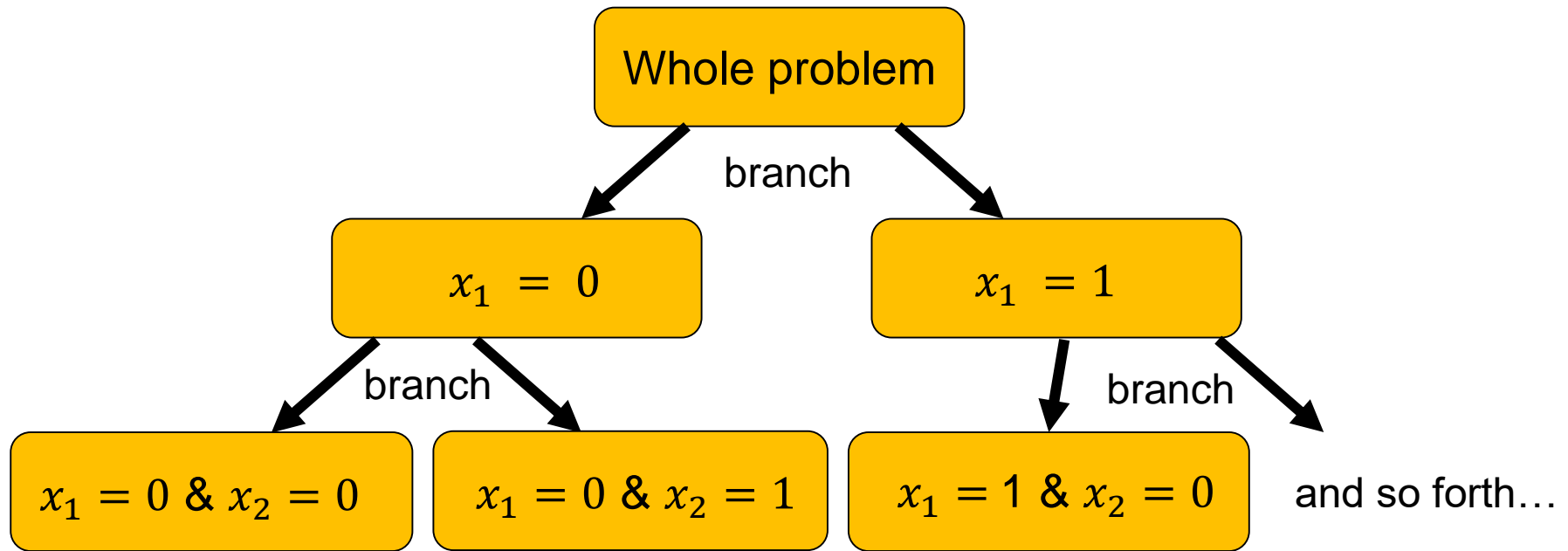
An Example: Branch&Bound for the KP

$$\max. \sum_{j=1}^n p_j x_j \text{ with } x_j \in \{0, 1\}$$

$$\text{s.t. } \sum_{j=1}^n w_j x_j \leq W$$

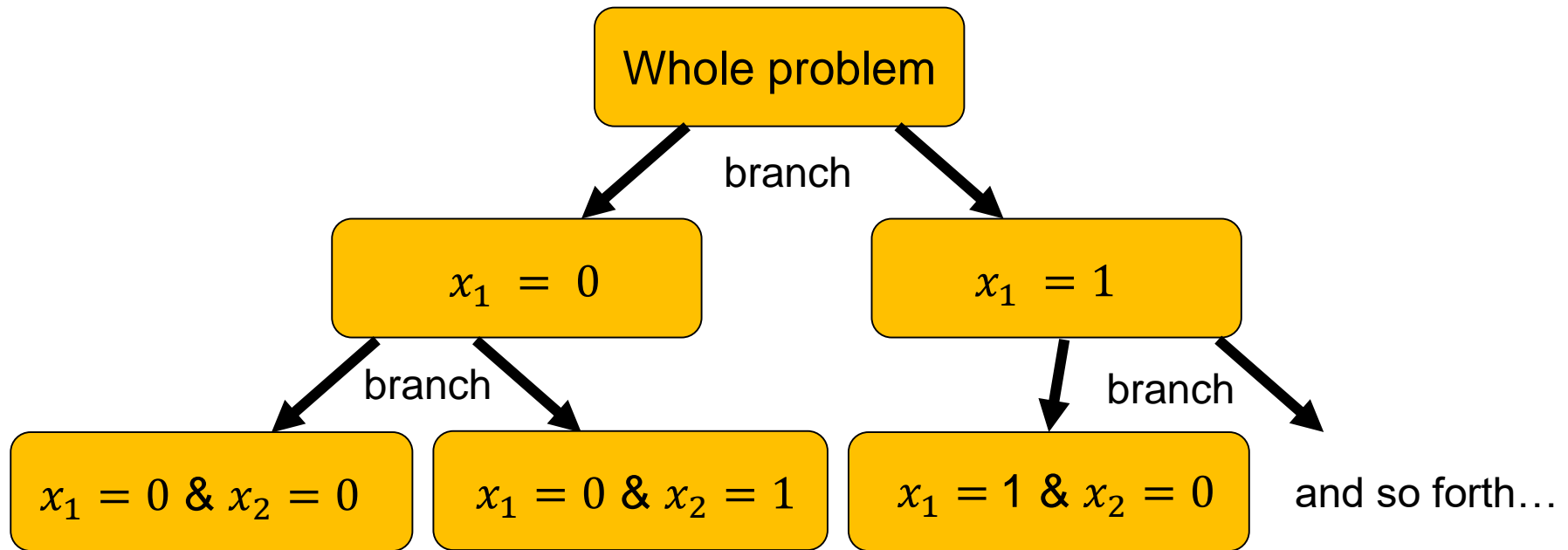


KP: How to Branch?



! order of variables plays an important role
optimally, the subproblems don't overlap

KP: How to Bound?



Maximization, so LB by greedy approach for example:

Choose items in decreasing profit/weight ratio until knapsack full

UB by relaxation of constraints (on the variables here):

Use greedy algorithm and pack add. item partially if there is space
...this variable can be used to branch next

Dynamic Programming

Dynamic Programming

Wikipedia:

“[...] **dynamic programming** is a method for solving a complex problem by breaking it down into a collection of simpler subproblems.”

But that's not all:

- dynamic programming also makes sure that the subproblems are not solved too often but only once by keeping the solutions of simpler subproblems in memory (“trading space vs. time”)
- it is an exact method, i.e. in comparison to the greedy approach, it always solves a problem to optimality

Two Properties Needed

Optimal Substructure

A solution can be constructed efficiently from optimal solutions of sub-problems

Overlapping Subproblems

Wikipedia: “[...] a problem is said to have **overlapping subproblems** if the problem can be broken down into subproblems which are reused several times or a recursive algorithm for the problem solves the same subproblem over and over rather than always generating new subproblems.”

Main Idea Behind Dynamic Programming

Main idea: solve larger subproblems by breaking them down to smaller, easier subproblems in a recursive manner

Typical Algorithm Design:

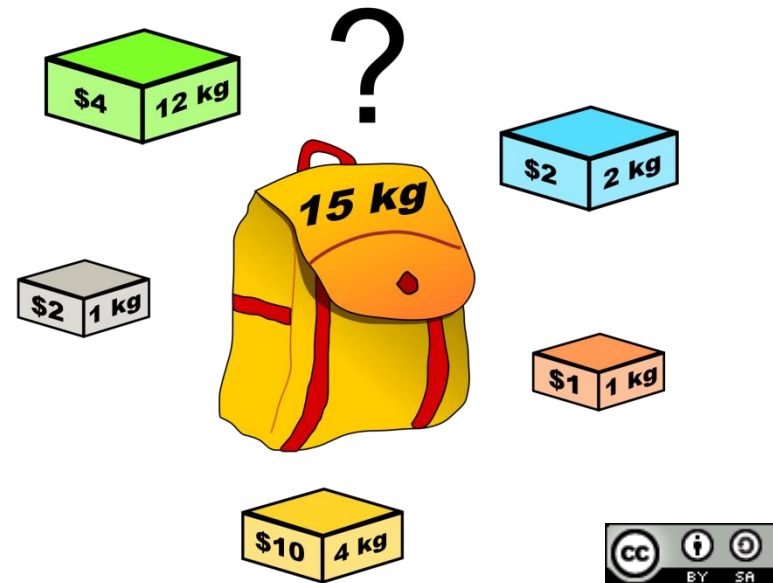
- ① decompose the problem into subproblems and think about how to solve a larger problem with the solutions of its subproblems
- ② specify how you compute the value of a larger problem recursively with the help of the optimal values of its subproblems (“Bellman equation”)
- ③ bottom-up solving of the subproblems (i.e. computing their optimal value), starting from the smallest by using the Bellman equality and a table structure to store the optimal values
- ④ eventually construct the final solution (can be omitted if only the value of an optimal solution is sought)

Example: The Knapsack Problem (KP)

Knapsack Problem

$$\max. \sum_{j=1}^n p_j x_j \text{ with } x_j \in \{0, 1\}$$

$$\text{s.t. } \sum_{j=1}^n w_j x_j \leq W$$



What are Good Subproblem Definitions for the KP?

Consider the following subproblems:

- 1) $P(i)$: optimal profit when packing exactly i items
- 2) $P(i)$: optimal profit when packing at most i items
- 3) $P(i, j)$: optimal profit when allowing to pack the first i items into a knapsack of size j

Which one allows us to solve larger subproblems from the solutions of smaller ones?

Which value are we actually interested in, when trying to solve the problem?

Opt. Substructure and Overlapping Subproblems

Consider the following subproblem:

$P(i, j)$: optimal profit when allowing to pack the first i items into a knapsack of size j

Optimal Substructure

The optimal choice of whether taking item i or not can be made easily for a knapsack of weight j if we know the optimal choice for items $1 \dots i - 1$:

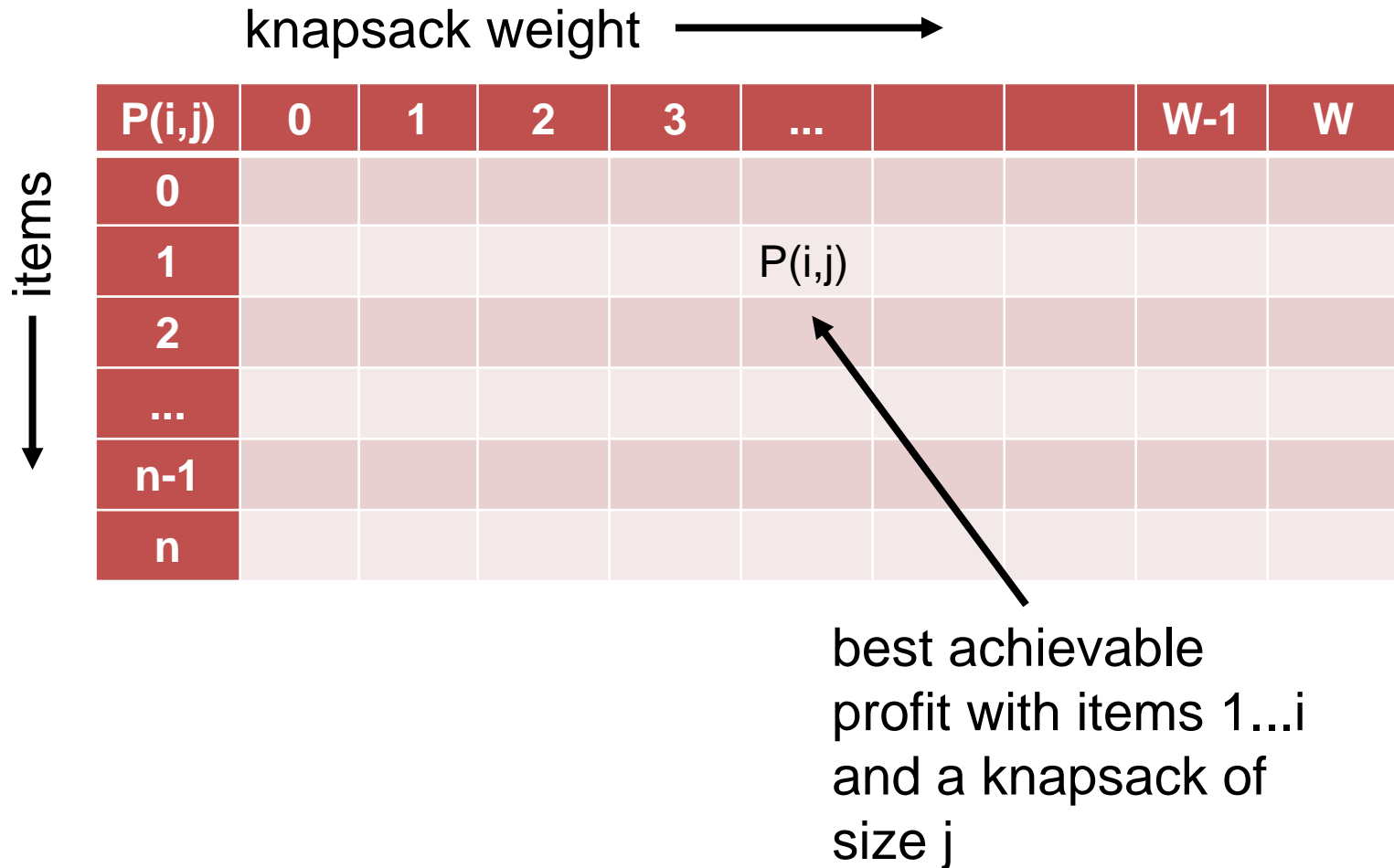
$$P(i, j) = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0 \\ P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Overlapping Subproblems

a recursive implementation of the Bellman equation is simple, but the $P(i, j)$ might need to be computed more than once!

Dynamic Programming Approach to the KP

To circumvent solving the subproblems more than once, we can store their results (in a matrix for example)...



Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W=11$.

knapsack weight \longrightarrow

$P(i,j)$	0	1	2	3	4	5	6	7	8	9	10	11
0												
1												
2												
3												
4												
5												

initialization:

$$P(i,j) = 0 \text{ if } i = 0 \text{ or } j = 0$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W=11$.

knapsack weight \longrightarrow

items \downarrow

$P(i,j)$	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0											
2	0											
3	0											
4	0											
5	0											

initialization:

$$P(i,j) = 0 \text{ if } i = 0 \text{ or } j = 0$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits
 (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

$P(i,j)$	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0										
2	0											
3	0											
4	0											
5	0											

items \downarrow

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

$P(i,j)$	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0										
2	0											
3	0											
4	0											
5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0									
2	0											
3	0											
4	0											
5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

$P(i,j)$	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0							
2	0											
3	0											
4	0											
5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4						
2	0											
3	0											
4	0											
5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4					
2	0											
3	0											
4	0											
5	0											

items \downarrow

A red arrow points from the cell (1,6) to (1,5) with the label $+p_1(=4)$. A blue arrow points from the cell (1,5) to (1,6).

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

	P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
items	0	0	0	0	0	0	0	0	0	0	0	0	0
	1	0	0	0	0	0	4	4	4	4	4	4	4
	2	0											
	3	0											
	4	0											
	5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4					
3	0											
4	0											
5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits
 (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4	10				
3	0											
4	0											
5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

	P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
items	0	0	0	0	0	0	0	0	0	0	0	0	0
	1	0	0	0	0	0	4	4	4	4	4	4	4
	2	0	0	0	0	0	4	4	10	10	10	10	10
	3	0											
	4	0											
	5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits
 (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

	P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	0	4	4	10	10	10	10	10
3	0	0	3	3	3								
4	0												
5	0												

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4	10	10	10	10	10
3	0	0	3	3	3	4						
4	0											
5	0											

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

items \downarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4	10	10	10	10	10
3	0	0	3	3	3	4	4					
4	0											
5	0											

Annotations: A red arrow points from the cell (3,5) to (3,6) with the label $+p_3 (= 3)$. A blue arrow points from the cell (2,6) to (3,6). The cell (3,6) is highlighted in green.

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4	10	10	10	10	10
3	0	0	3	3	3	4	4	10	etc.			
4	0											
5	0											

items \downarrow

Annotations: A red arrow points from the cell (3,6) to (3,7) with the label $+p_3 (= 3)$. A blue arrow points from the cell (3,7) to (2,7).

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits (5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

items \downarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4	10	10	10	10	10
3	0	0	3	3	3	4	4	10	10	13	13	13
4	0	0	3	3	5	5	8	10	10	13	13	15
5	0	0	3	3	5	6	8	10	10	13	13	15

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Example instance with 5 items with weights and profits
(5,4), (7,10), (2,3), (4,5), and (3,3). Weight restriction is $W = 11$.

knapsack weight \longrightarrow

items \downarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4	10	10	10	10	10
3	0	0	3	3	3	4	4	10	10	13	13	13
4	0	0	3	3	5	5	8	10	10	13	13	15
5	0	0	3	3	5	6	8	10	10	13	13	15

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

Dynamic Programming Approach to the KP

Question: How to obtain the actual packing?

Answer: we just need to remember where the max came from!

knapsack weight \longrightarrow

P(i,j)	0	1	2	3	4	5	6	7	8	9	10	11
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	4	4	4	4	4	4	4
2	0	0	0	0	0	4	4	10	10	10	10	10
3	0	0	3	3	3	4	4	10	10	13	13	13
4	0	0	3	3	5	5	8	10	10	13	13	15
5	0	0	3	3	5	6	8	10	10	13	13	15

items \downarrow

Annotations: $x_1 = 0$ (blue arrow from (1,0) to (1,1)), $x_2 = 1$ (red arrow from (1,6) to (2,6)), $x_3 = 0$ (blue arrow from (2,7) to (2,8)), $x_4 = 1$ (red arrow from (2,10) to (3,10)), $x_5 = 0$ (blue arrow from (3,11) to (3,12)).

for $i = 1$ to n :

for $j = 1$ to W :

$$P(i, j) = \begin{cases} P(i - 1, j) & \text{if } w_i > j \\ \max\{P(i - 1, j), p_i + P(i - 1, j - w_i)\} & \text{if } w_i \leq j \end{cases}$$

(Randomized) Search Heuristics

Slides with this light blue background have not been discussed in the lecture and are thus not part of the exam.

I left them in for those of you who are interested to learn about the subject anyway.

Motivation General Search Heuristics

- often, problem complicated and not much time available to develop a problem-specific algorithm
- search heuristics are a good choice:
 - relatively **easy to implement**
 - **easy to adapt/change/improve**
 - e.g. when the problem formulation changes in an early product design phase
 - or when slightly different problems need to be solved over time
- randomized/stochastic algorithms are a good choice because they are robust to noise

Which algorithms will we touch?

- ➊ Randomized Local Search (RLS)
- ➋ Variable Neighborhood Search (VNS)
- ➌ Tabu Search (TS)
- ➍ Evolutionary Algorithms (EAs)

Neighborhoods

For most (stochastic) search heuristics, we need to define a *neighborhood structure*

- which search points are close to each other?

Example: k-bit flip / Hamming distance k neighborhood

- search space: bitstrings of length n ($\Omega = \{0,1\}^n$)
- two search points are neighbors if their **Hamming distance** is k
- in other words: x and y are neighbors if we can flip exactly k bits in x to obtain y
- 0001001101 is neighbor of
0001000101 for $k=1$
0101000101 for $k=2$
1101000101 for $k=3$

Neighborhoods II

Example: possible neighborhoods for the **knapsack problem**

- search space again bitstrings of length n ($\Omega = \{0,1\}^n$)
- **Hamming distance 1 neighborhood:**
 - add an item or remove it from the packing
- **replacing 2 items neighborhood:**
 - replace one chosen item with an unchosen one
 - makes only sense in combination with other neighborhoods because the number of items stays constant
- **Hamming distance 2 neighborhood** on the contrary:
 - allows to change 2 arbitrary items, e.g.
 - add 2 new items
 - remove 2 chosen items
 - or replace one chosen item with an unchosen one

Randomized Local Search (RLS)

Idea behind (Randomized) Local Search:

- explore the local neighborhood of the current solution (randomly)

Pure Random Search:

- go to randomly chosen neighbor

First Improvement Local Search:

- go to first (randomly) chosen neighbor which is better

Best Improvement strategy:

- always go to the best neighbor
- not random anymore
- computationally expensive if neighborhood large

Variable Neighborhood Search

Main Idea: [N. Mladenovic and P. Hansen, 1997]

- change the neighborhood from time to time
 - local optima not necessarily the same for different neighborhood operators
 - but often close to each other
 - global optimum is local optimum for all neighborhoods
- rather a framework than a concrete algorithm
 - e.g. deterministic and stochastic neighborhood changes
- typically combined with (i) first improvement, (ii) a random order in which the neighbors are visited and (iii) restarts

N. Mladenovic and P. Hansen (1997). "Variable neighborhood search". *Computers and Operations Research* 24 (11): 1097–1100.

Disadvantages of local searches (with or without varying neighborhoods)

- they get stuck in local optima
- have problems to traverse large plateaus of equal objective function value (“random walk”)

Tabu search addresses these by

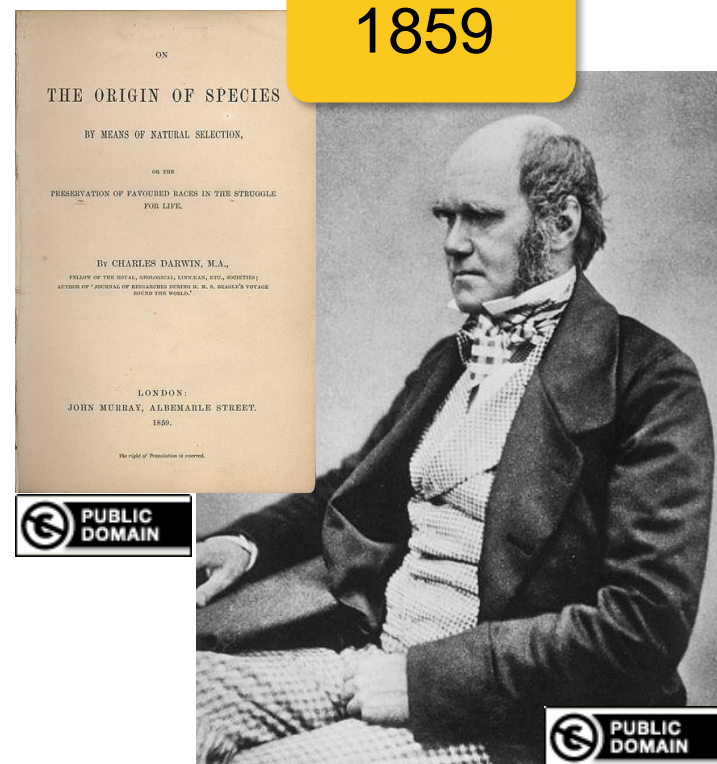
- allowing worsening moves if all neighbors are explored
- introducing a tabu list of temporarily not allowed moves
- those restricted moves are
 - problem-specific and
 - can be specific solutions or not permitted “search directions” such as “don’t include this edge anymore” or “do not flip this specific bit”
- the tabu list is typically restricted in size and after a while, restricted moves are permitted again

Stochastic Optimization Algorithms

One class of (bio-inspired) stochastic optimization algorithms: Evolutionary Algorithms (EAs)

- Class of optimization algorithms originally inspired by the idea of **biological evolution**
- selection, mutation, recombination

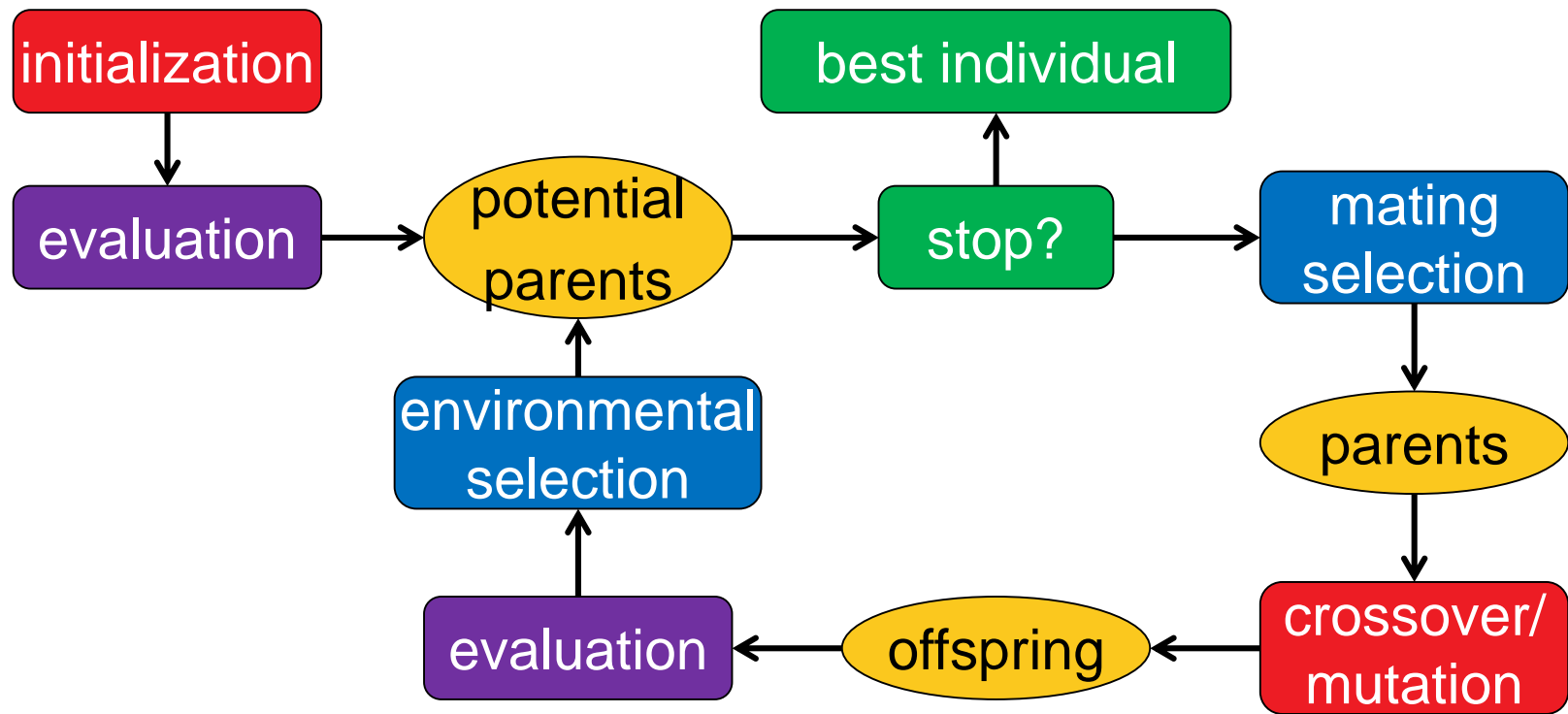
1859



Metaphors

Classical Optimization	Evolutionary Computation
variables or parameters	variables or chromosomes
candidate solution vector of decision variables / design variables / object variables	individual, offspring, parent
set of candidate solutions	population
objective function loss function cost function error function	fitness function
iteration	generation

Generic Framework of an EA



stochastic operators

“Darwinism”

stopping criteria

Important:
representation (search space)

The Historic Roots of EAs

Genetic Algorithms (GA)

J. Holland 1975 and D. Goldberg (USA)

$$\Omega = \{0, 1\}^n$$

Evolution Strategies (ES)

I. Rechenberg and H.P. Schwefel, 1965 (Berlin)

$$\Omega = \mathbb{R}^n$$

Evolutionary Programming (EP)

L.J. Fogel 1966 (USA)

Genetic Programming (GP)

J. Koza 1990 (USA)

$\Omega = \text{space of all programs}$

nowadays one umbrella term: **evolutionary algorithms**

Note: Handling Constraints

Several generic ways to handle constraints, e.g.:

- **resampling** until a new feasible point is found (“often bad idea”)
- **penalty function** approach: add constraint violation term (potentially scaled)
- **repair** approach: after generation of a new point, repair it (e.g. with a heuristic) to become feasible again if infeasible
 - continue to use repaired solution in the population or
 - use repaired solution only for the evaluation?
- **multiobjective** approach: keep objective function and constraint functions separate and try to optimize all of them in parallel
- ...

Examples for some EA parts

Selection

Selection is the major determinant for specifying the trade-off between **exploitation** and **exploration**

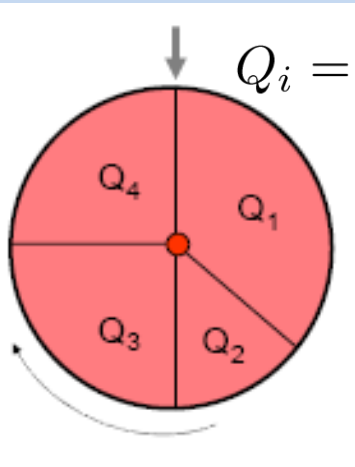
Selection is either

stochastic

or

deterministic

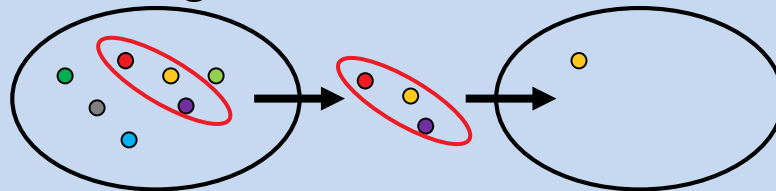
e.g. fitness proportional



$$Q_i = \frac{f(x_i)}{\sum_{j=1}^{\mu} f(x_j)}$$

Disadvantage:
depends on
scaling of f

e.g. via a tournament



e.g. $(\mu+\lambda)$, (μ, λ)



Mating selection (selection for variation): usually stochastic

Environmental selection (selection for survival): often deterministic

Variation Operators

Variation aims at generating new individuals on the basis of those individuals selected for mating

Variation = Mutation and Recombination/Crossover

mutation: $mut: \Omega \rightarrow \Omega$

recombination: $recomb: \Omega^r \rightarrow \Omega^s$ where $r \geq 2$ and $s \geq 1$

- choice always depends on the problem and the chosen representation
- however, there are some operators that are applicable to a wide range of problems and tailored to **standard representations** such as vectors, permutations, trees, etc.

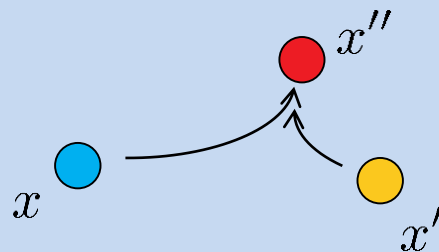
Variation Operators: Guidelines

Two desirable properties for **mutation** operators:

- every solution can be generation from every other with a probability greater than 0 (“exhaustiveness”)
- $d(x, x') < d(x, x'') \Rightarrow \text{Prob}(\text{mut}(x) = x') > \text{Prob}(\text{mut}(x) = x'')$ (“locality”)

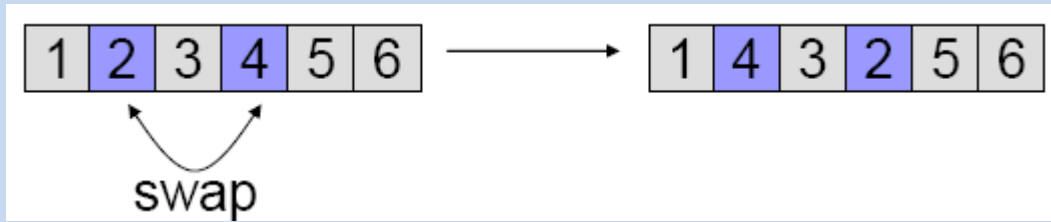
Desirable property of **recombination** operators (“in-between-ness”):

$$x'' = \text{recomb}(x, x') \Rightarrow d(x'', x) \leq d(x, x') \wedge d(x'', x') \leq d(x, x')$$

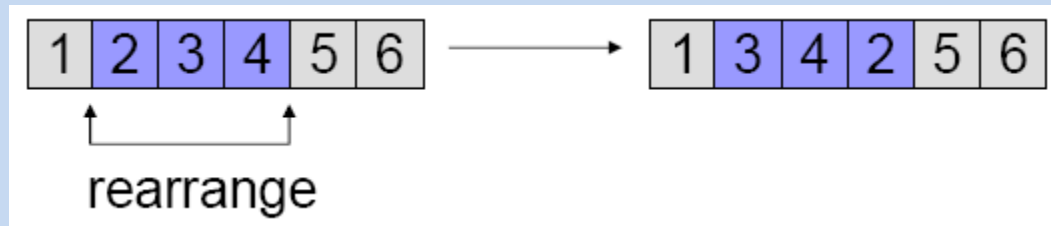


Examples of Mutation Operators on Permutations

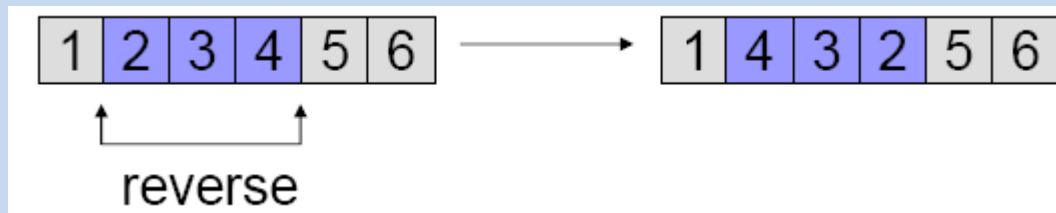
Swap:



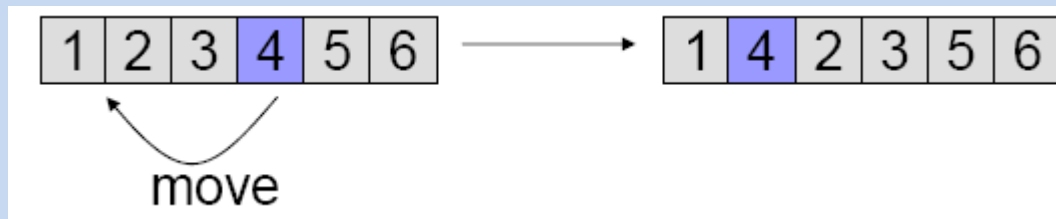
Scramble:



Invert:

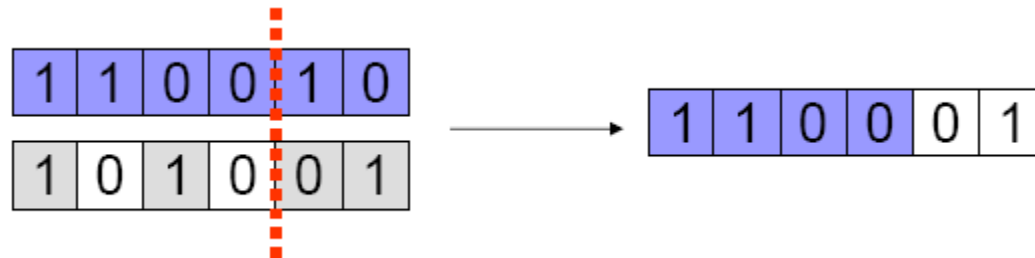


Insert:

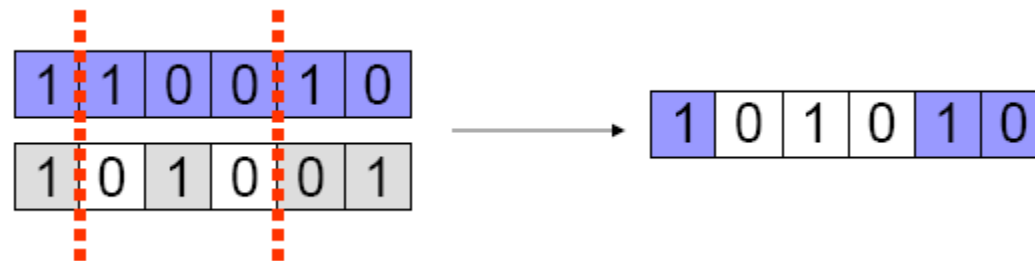


Examples of Recombination Operators: $\{0,1\}^n$

1-point crossover



n-point crossover



uniform crossover



choose each bit independently from one parent or another

A Canonical Genetic Algorithm

- binary search space, maximization
- uniform initialization
- generational cycle: of the population
 - evaluation of solutions
 - mating selection (e.g. roulette wheel)
 - crossover (e.g. 1-point)
 - environmental selection (e.g. plus-selection)

Full Circle: CMA-ES to solve Continuous Problems

A stochastic blackbox search template to minimize $f: \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize distribution parameters θ , set population size $\lambda \in \mathbb{N}$

While happy do:

- Sample distribution $P(\mathbf{x}|\theta) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda \in \mathbb{R}^n$
- Evaluate $\mathbf{x}_1, \dots, \mathbf{x}_\lambda$ on f
- Update parameters $\theta \leftarrow F_\theta(\theta, \mathbf{x}_1, \dots, \mathbf{x}_\lambda, f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda))$

For CMA-ES and evolution strategies in general:

sample distributions = multivariate Gaussian distributions

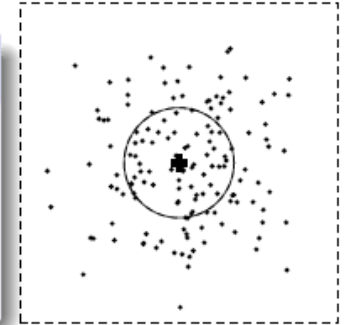
Sampling New Candidate Solutions (Offspring)

Evolution Strategies

New search points are sampled normally distributed

$$\mathbf{x}_i \sim \mathbf{m} + \sigma \mathcal{N}_i(\mathbf{0}, \mathbf{C}) \quad \text{for } i = 1, \dots, \lambda$$

as perturbations of \mathbf{m} , where $\mathbf{x}_i, \mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, $\mathbf{C} \in \mathbb{R}^{n \times n}$



where

- the **mean** vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution
- the so-called **step-size** $\sigma \in \mathbb{R}_+$ controls the *step length*
- the **covariance matrix** $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the **shape** of the distribution ellipsoid

here, all new points are sampled with the same parameters

from [Auger, p. 10]

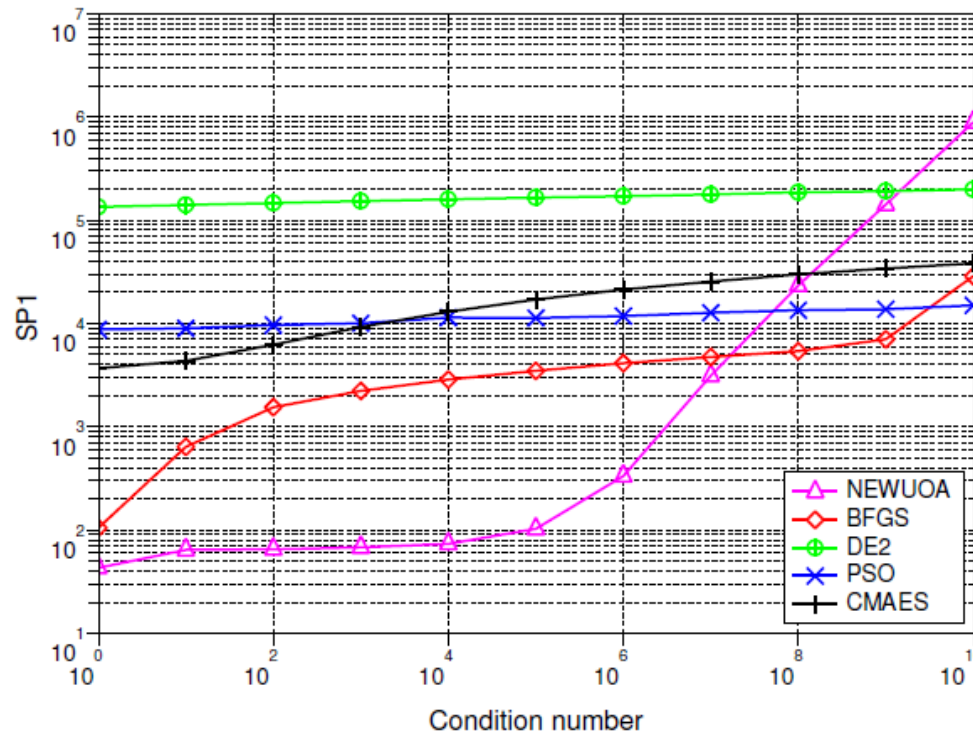
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

f convex quadratic, separable with varying condition number α

Ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

$f(x) = g(x^T H x)$ with

H diagonal

g identity (for **BFGS** and **NEWUOA**)

g any order-preserving = strictly increasing function (for all other)

SP1 = average number of objective function evaluations¹⁴ to reach the target function value of $g^{-1}(10^{-9})$

from [Nikolaus Hansen]

¹⁴ Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

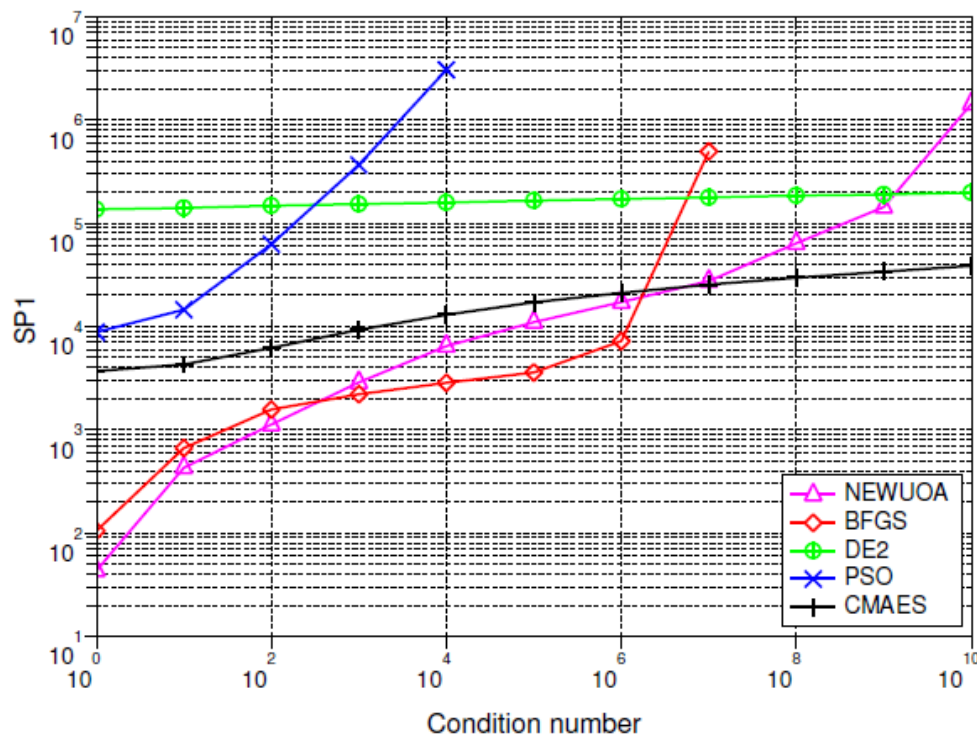
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

f convex quadratic, non-separable (rotated) with varying condition number α

Rotated Ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

$f(x) = g(x^T H x)$ with

H full

g identity (for **BFGS** and **NEWUOA**)

g any order-preserving = strictly increasing function (for all other)

SP1 = average number of objective function evaluations¹⁵ to reach the target function value of $g^{-1}(10^{-9})$

from [Nikolaus Hansen]

¹⁵ Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

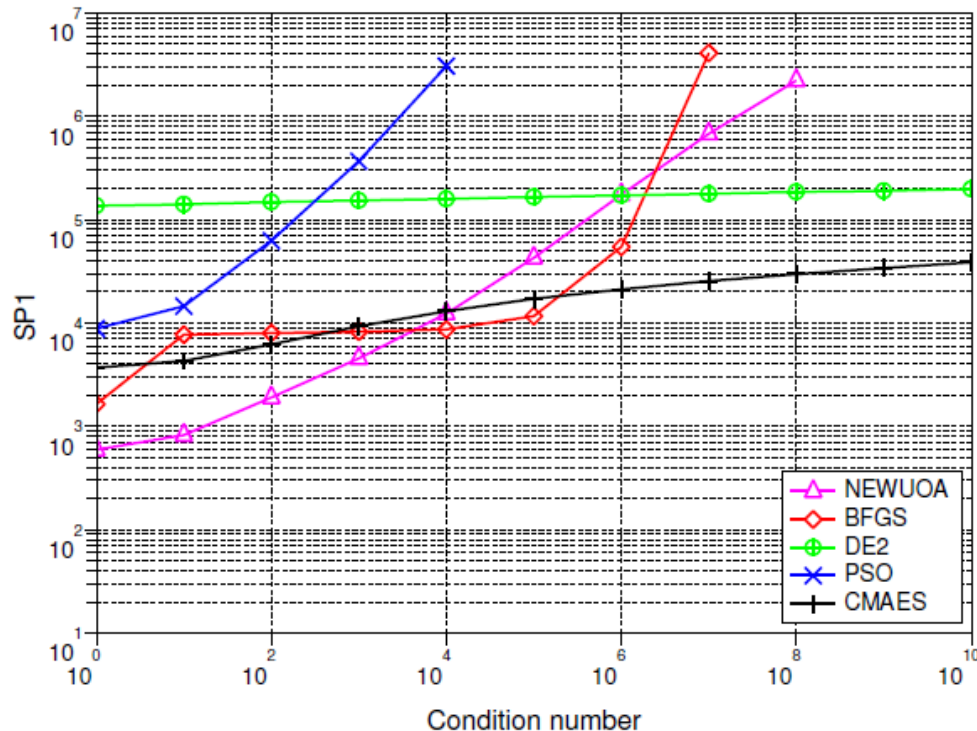
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

f non-convex, non-separable (rotated) with varying condition number α

Sqrt of sqrt of rotated ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

$f(x) = g(x^T H x)$ with

H full

$g : x \mapsto x^{1/4}$ (for **BFGS** and

NEWUOA)

g any order-preserving = strictly increasing function (for all other)

SP1 = average number of objective function evaluations¹⁶ to reach the target function value of $g^{-1}(10^{-9})$

from [Nikolaus Hansen]

¹⁶Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

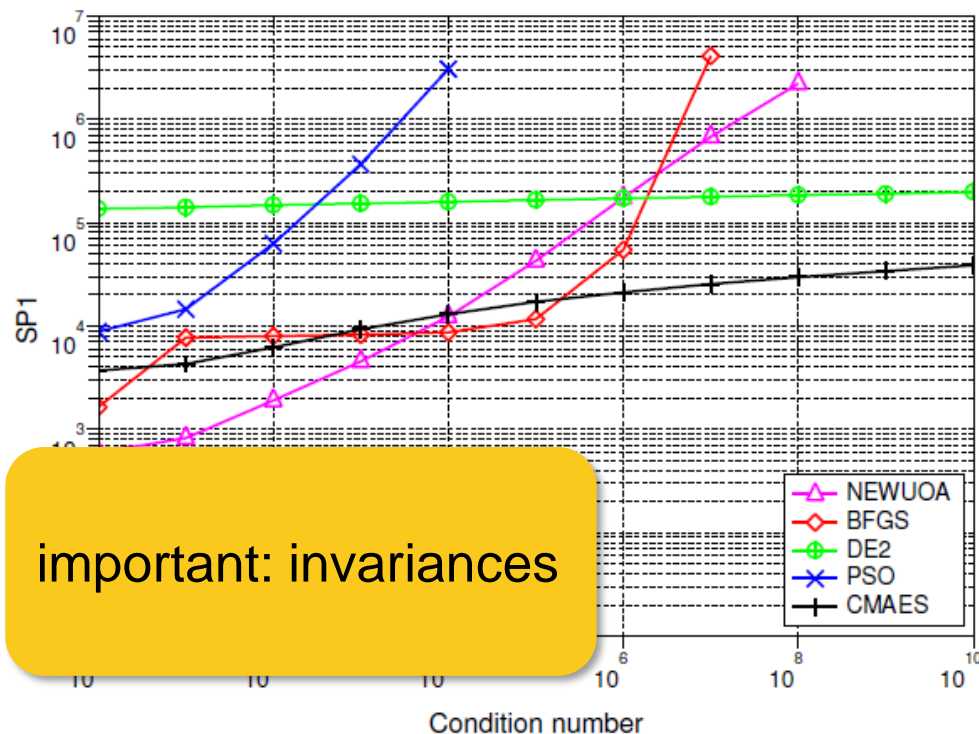
Influence of Condition Number + Invariance

Comparing Experiments

Comparison to BFGS, NEWUOA, PSO and DE

f non-convex, non-separable (rotated) with varying condition number α

Sqrt of sqrt of rotated ellipsoid dimension 20, 21 trials, tolerance $1e-09$, eval max $1e+07$



important: invariances

BFGS (Broyden et al 1970)

NEWUOA (Powell 2004)

DE (Storn & Price 1996)

PSO (Kennedy & Eberhart 1995)

CMA-ES (Hansen & Ostermeier 2001)

$f(x) = g(x^T H x)$ with

H full

$g : x \mapsto x^{1/4}$ (for **BFGS** and

NEWUOA)

g any order-preserving = strictly increasing function (for all other)

SP1 = average number of objective function evaluations¹⁶ to reach the target function value of $g^{-1}(10^{-9})$

from [Nikolaus Hansen]

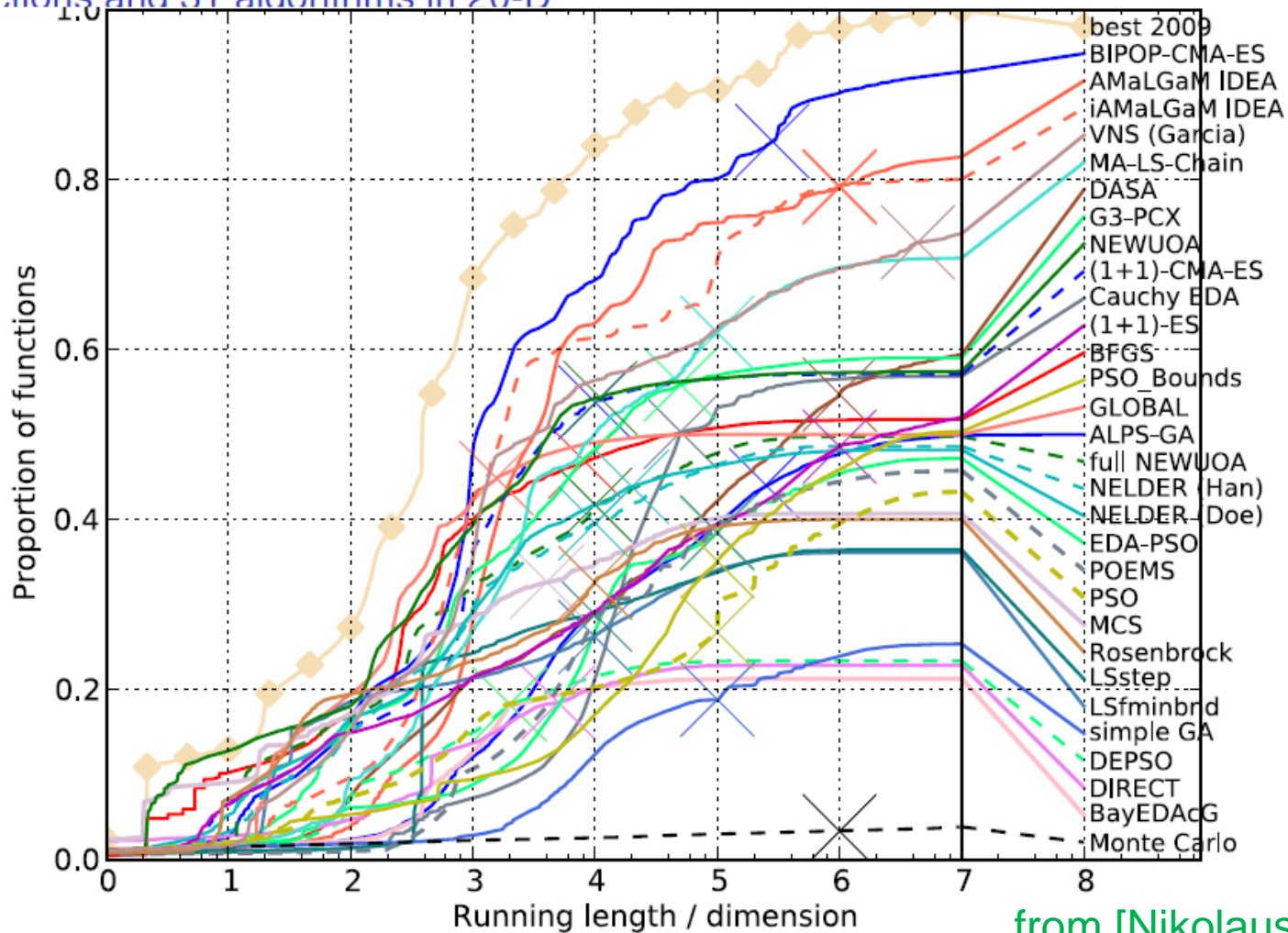
¹⁶ Auger et.al. (2009): Experimental comparisons of derivative free optimization algorithms, SEA

Performance on BBOB Testbed: Data Profile

Comparing Experiments

Comparison during BBOB at GECCO 2009

24 functions and 31 algorithms in 20-D



from [Nikolaus Hansen]