

Derivative Free Optimization

**Optimization Master Paris Saclay / IPP
AMS Master
2023/2024**

Anne Auger
RandOpt team
Inria and CMAP, Ecole Polytechnique, IP Paris
anne.auger@inria.fr

Organization of the class

When: Friday afternoon - 2pm - 5:15pm at ENSTA

01/12/2023	room 1320
08/12/2023	room 1321
15/12/2023	room 1320
22/12/2022	room 1320
12/01/2024	room 1320
19/01/2024	room 1320
26/01/2024	room 1320
02/02/2024	room 1320
09/02/2024	room 1320
16/02/2024 [EXAM]	

Evaluation

Written exam on 16/02/2024

Project (in group) around benchmarking/testing of algorithms

- oral presentation to the class

Topics covered

Derivative Free Optimization / Black-box optimization

Single-objective optimization

what makes a problem difficult

algorithm to solve those difficulties (mostly stochastic)

Multi-objective optimization [taught D. Brockhoff]

Benchmarking (partly taught by D. Brockhoff)

Practical Exercises [bring your laptops]

practical exercises: **implement/manipulate** algorithms

Python / Matlab / ...

ultimate goal: optimize a (real) black-box problem on your own

- understand and visualize convergence / adaptation / invariance
- experience numerics numerical errors, finite machine precision

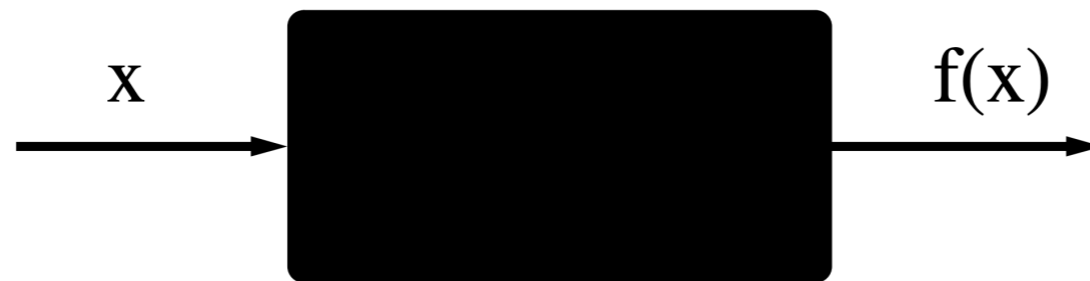
Derivative-Free / Black-box Optimization

Task: minimize a numerical **objective** function (also called *fitness* function or *loss* function)

$$f: \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}, x \mapsto f(x) \in \mathbb{R}$$

without derivatives (gradient). Ω : search space, n : dimension of the search space

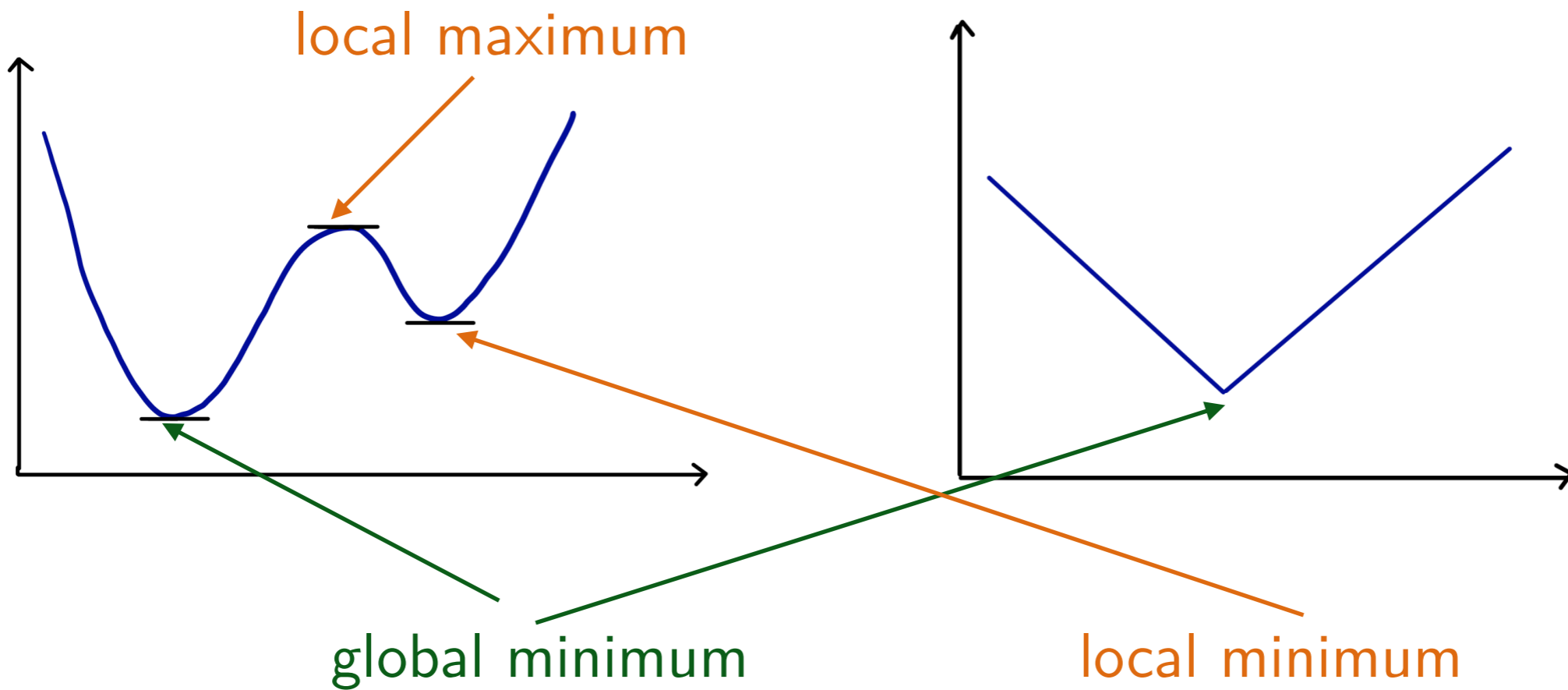
Also called **zero-order black-box** optimization



The function is seen by the algorithm as a zero-order **oracle** [a first order oracle would also return gradients] that can be queried at points and the oracle returns an answer

Reminder: Local versus Global Optimum

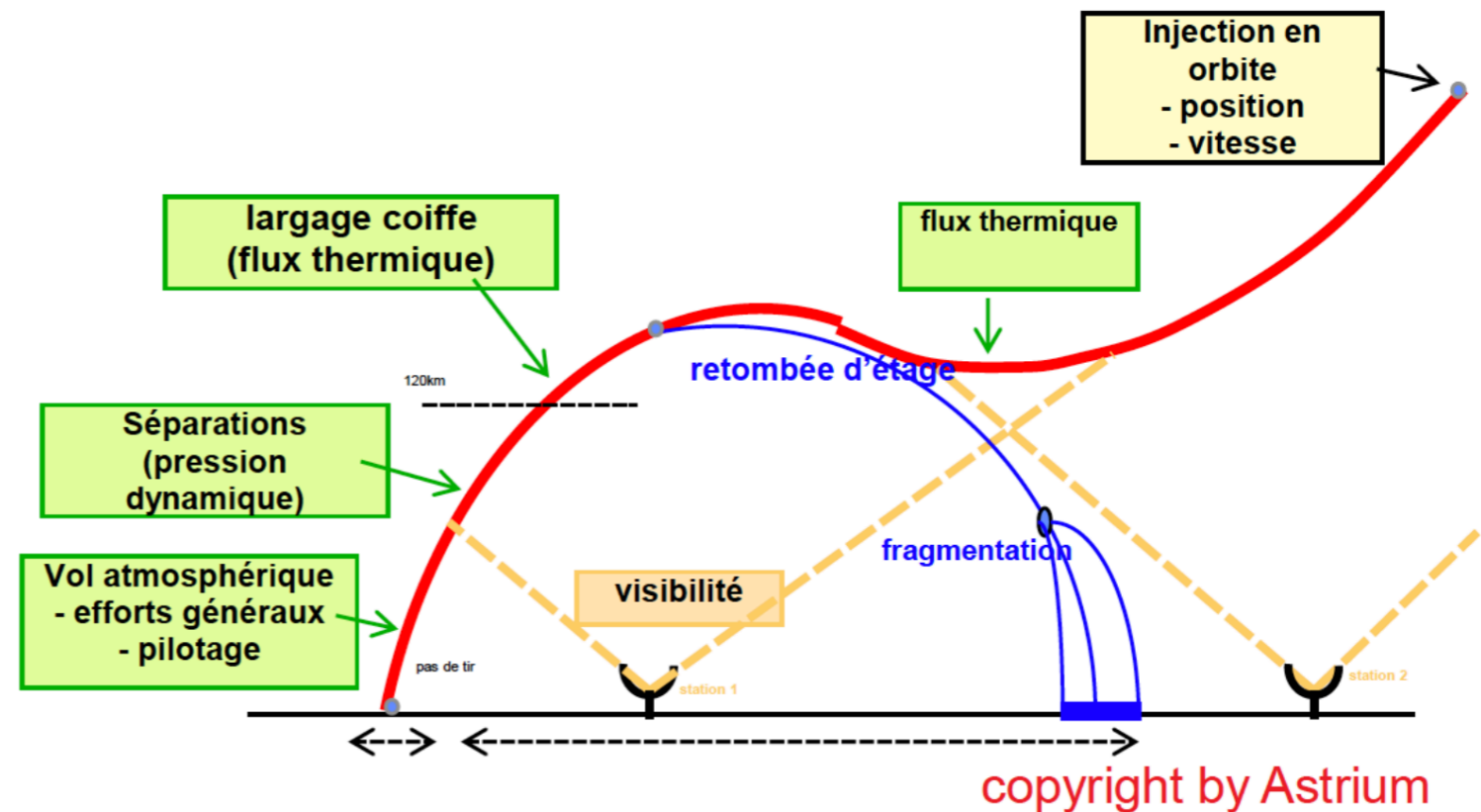
$n=1$



Examples: Optimization of the Design of a Launcher



Poppy

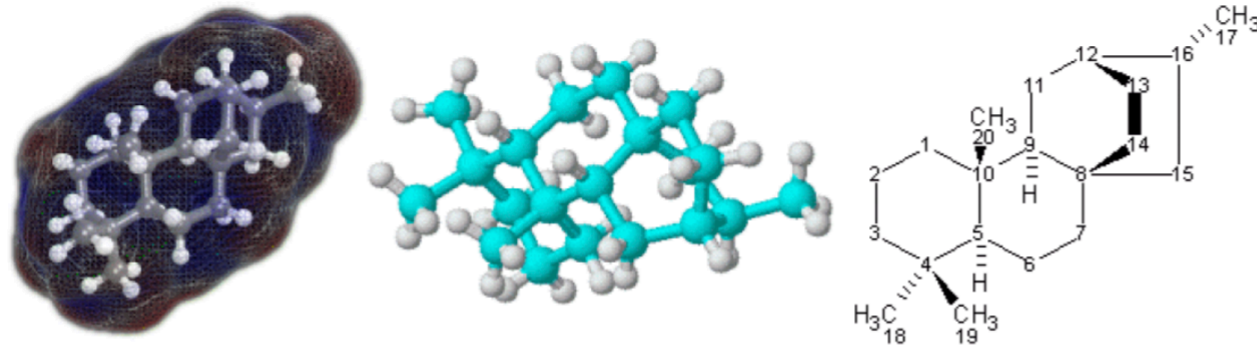


- Scenario: multi-stage launcher brings a satellite into orbit
- Minimize the overall cost of a launch
- Parameters: propellant mass of each stage / diameter of each stage / flux of each engine / parameters of the command law

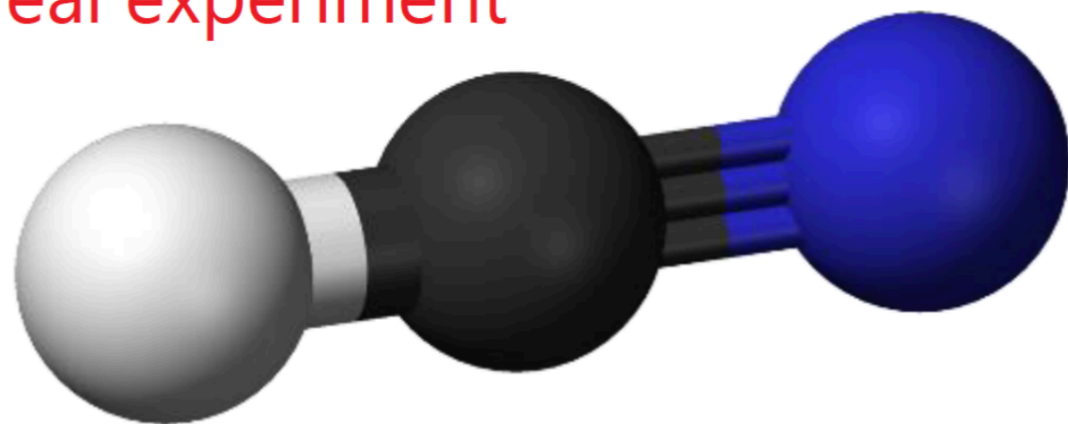
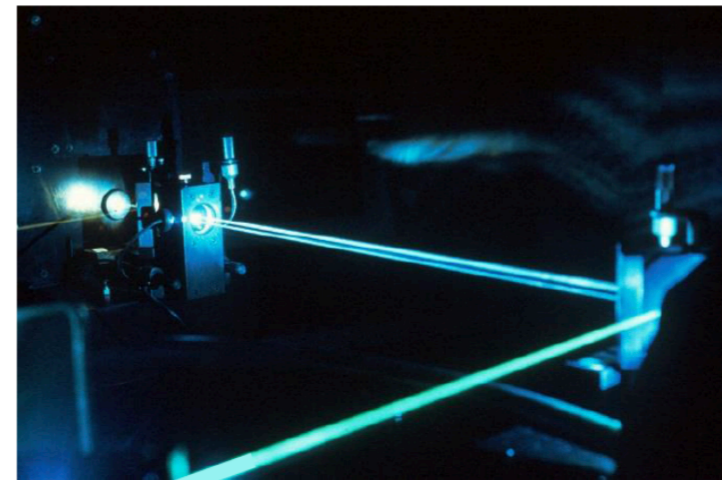
*23 continuous parameters to optimize
+ constraints*

Control of the Alignment of Molecules

application domain: quantum physics or chemistry



Objective function:
via **numerical simulation**
or a **real experiment**



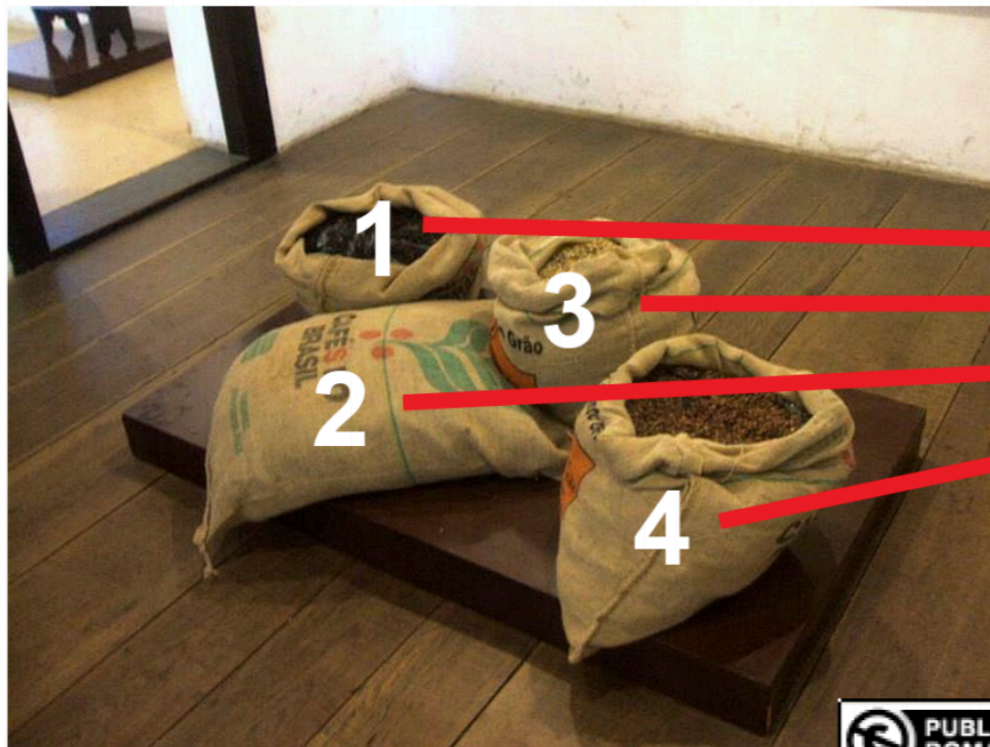
possible application in drug design

*In the case of a real lab experiment: the objective function is
a real black-box*

Coffee Tasting Problem (A real Black-box)

Coffee Tasting Problem

- ▶ Find a mixture of coffee in order to keep the coffee taste from one year to another
- ▶ Objective function = opinion of one expert



 PUBLIC DOMAIN



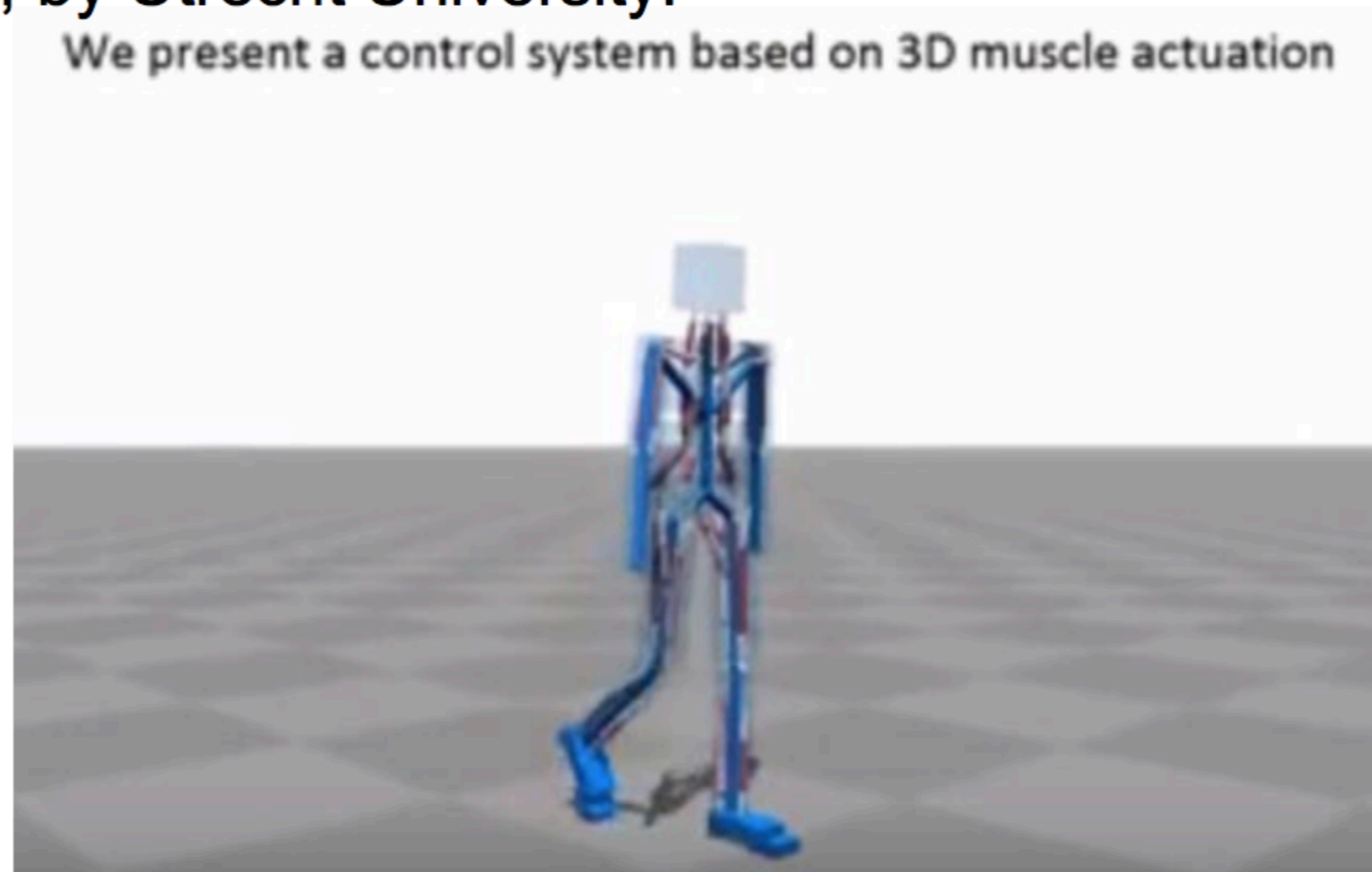
 CC BY SA

Quasipalm

M. Herdy: "Evolution Strategies with subjective selection", 1996

A last Application

Computer simulation teaches itself to walk upright (virtual robots (of different shapes) learning to walk, through stochastic optimization (CMA-ES)), by Utrecht University:



<https://www.youtube.com/watch?v=yci5Fu11ovk>

T. Geitjenbeek, M. Van de Panne, F. Van der Stappen: "Flexible Muscle-Based Locomotion for Bipedal Creatures", SIGGRAPH Asia, 2013.

What is the Goal?

- We want to find x^\star such that $f(x^\star) \leq f(x)$ for all x

$$x^\star \in \operatorname{argmin}_x f(x)$$

- In general we will never find x^\star

why?

What is the Goal?

- We want to find x^\star such that $f(x^\star) \leq f(x)$ for all x

$$x^\star \in \operatorname{argmin}_x f(x)$$

- In general we will never find x^\star
- Because of the numerical/continuous nature of the search space we typically never hit exactly x^\star , we instead converge to a solution:

we want to find $x_t \in \mathbb{R}^n$ such that $\lim_{t \rightarrow \infty} f(x_t) = \min f$

of course we want **fast** convergence

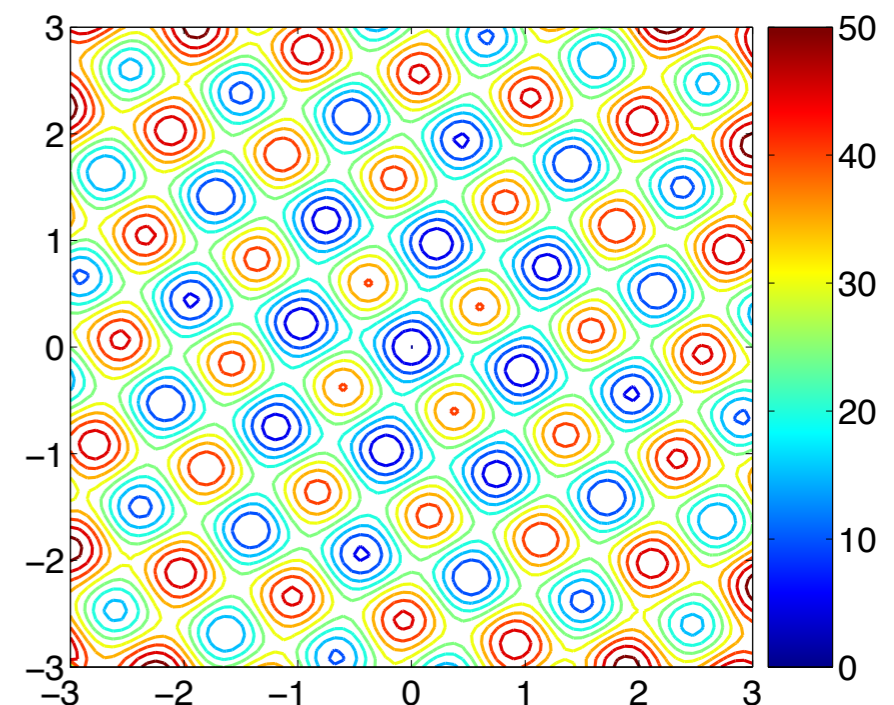
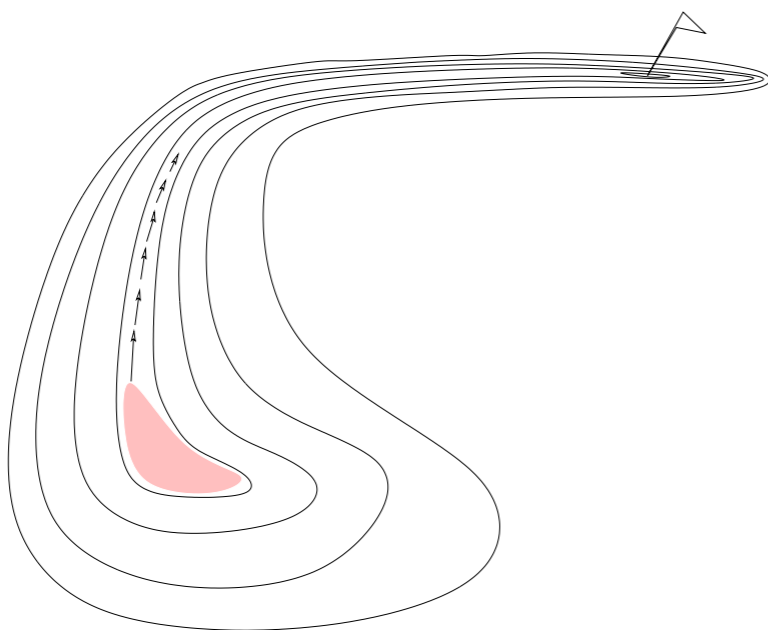
Level Sets of a Function

Level Sets: Visualization of a Function

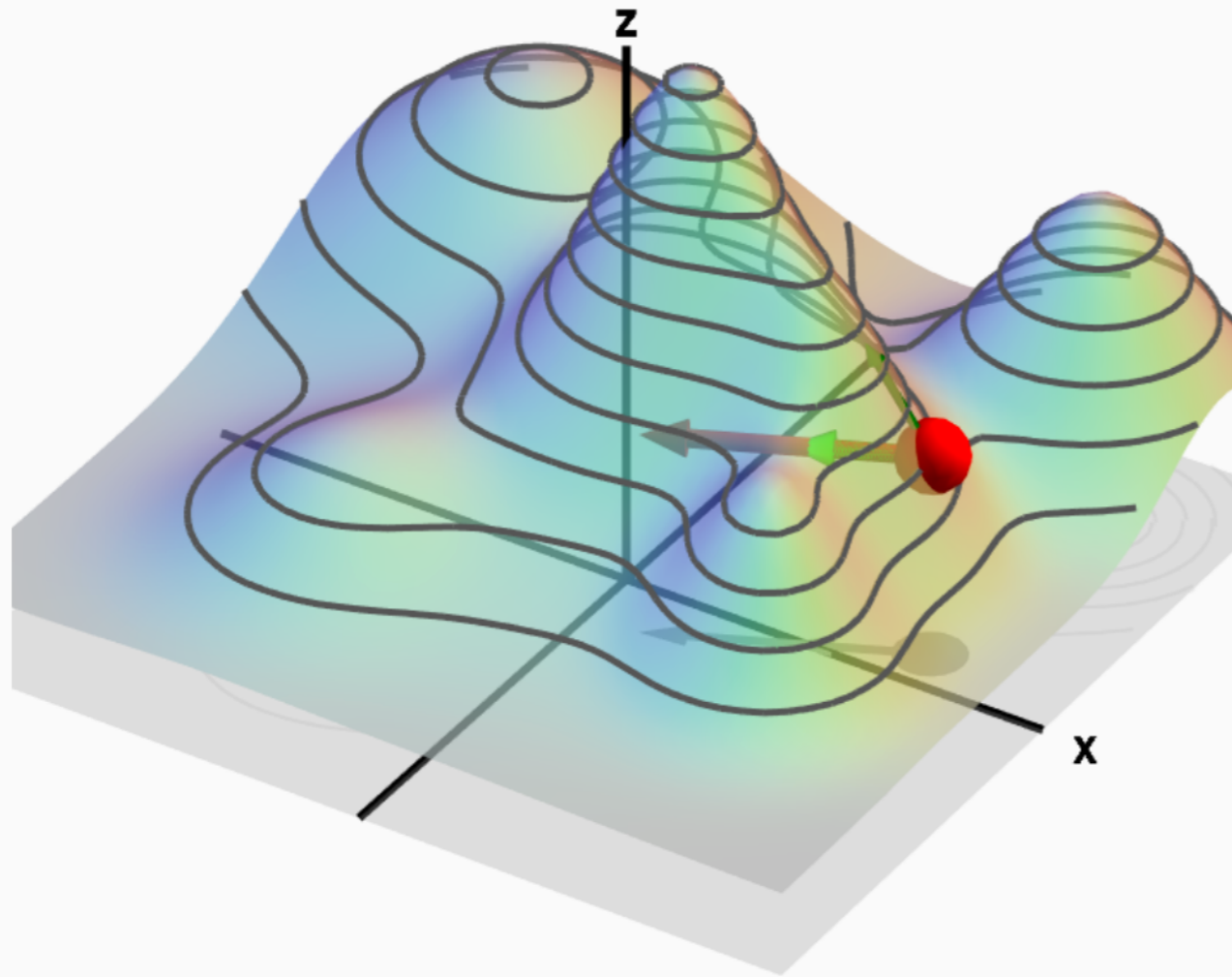
One-dimensional (1-D) representations are often misleading (as 1-D optimization is “trivial”, see slides related to curse of dimensionality), we therefore often represent **level-sets** of functions

$$\mathcal{L}_c = \{x \in \mathbb{R}^n \mid f(x) = c, \}, c \in \mathbb{R}$$

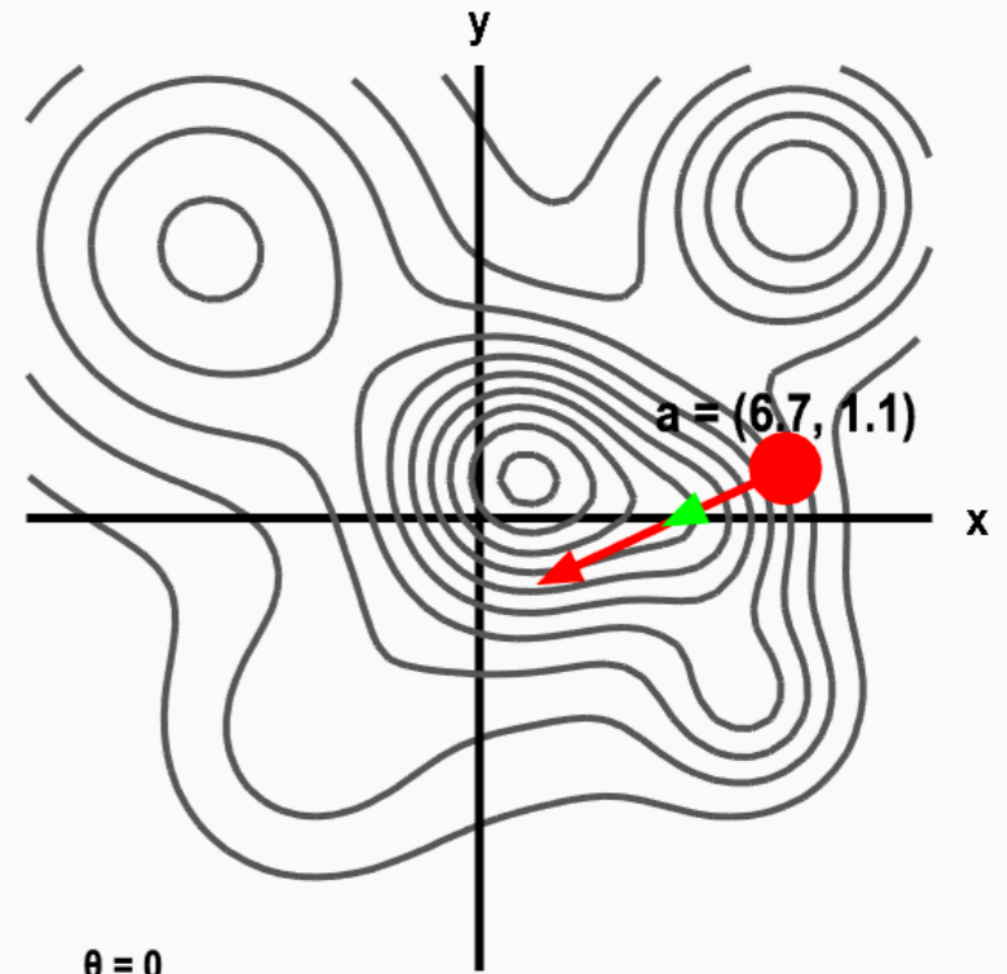
Examples of level sets in 2D



Level Sets: Visualization of a Function



$\theta = 0$
 $u = (-0.91, -0.42)$
 $a = (6.7, 1.1)$
 $\nabla f(a) = (-1.81, -0.85)$
 $D_u f(a) = 2.00$
 $|\nabla f(a)| = 2.00$



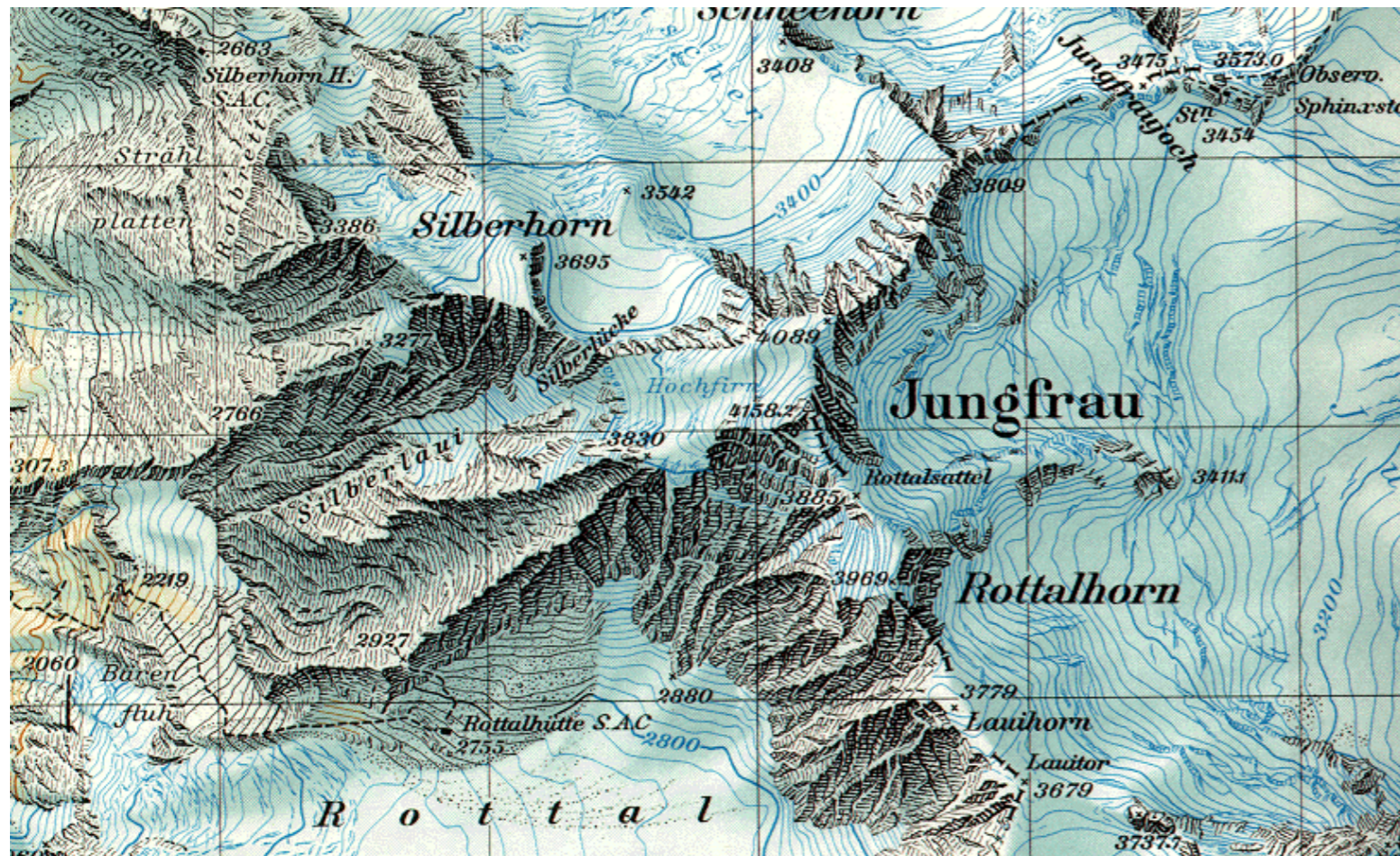
$\theta = 0$
 $u = (-0.91, -0.42)$
 $\nabla f(a) = (-1.81, -0.85)$
 $D_u f(a) = 2.00$
 $|\nabla f(a)| = 2.00$
 $f(a) = 4.87$



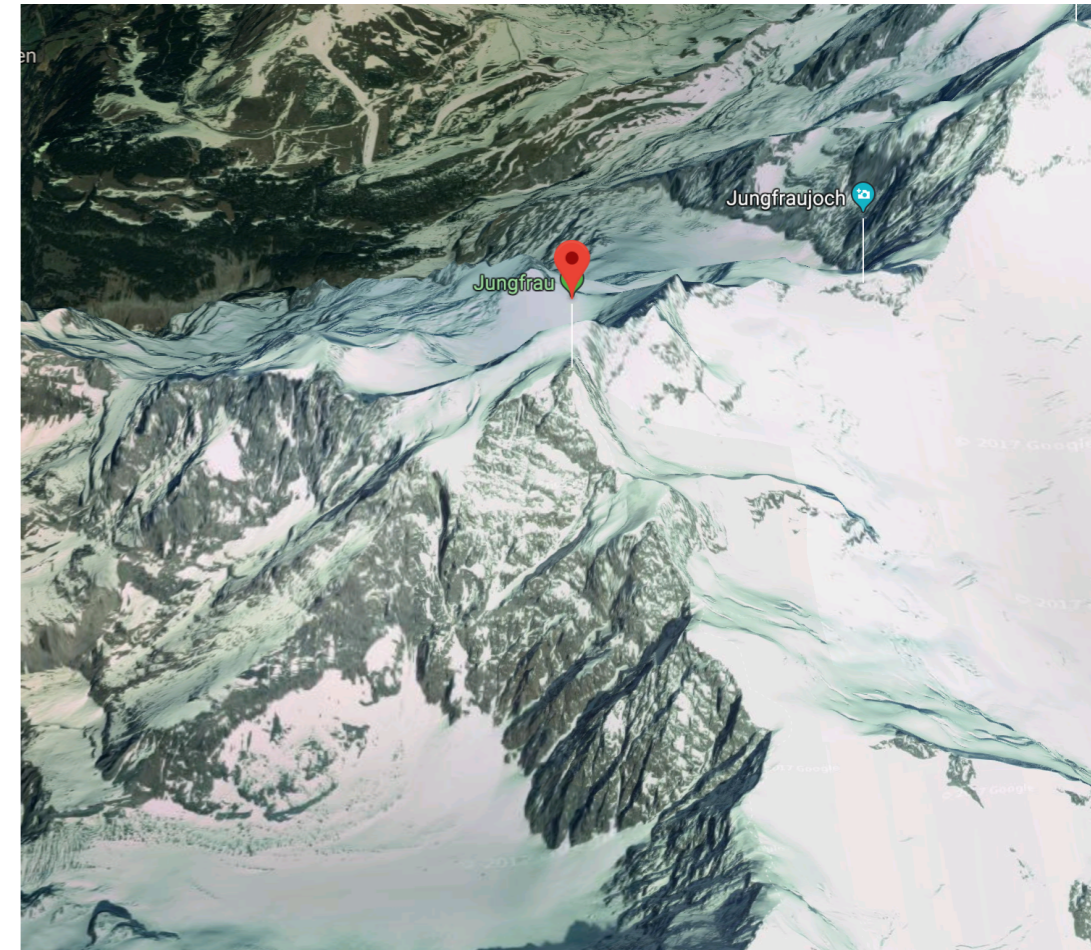
Source: Nykamp DQ, "Directional derivative on a mountain." From *Math Insight*. http://mathinsight.org/applet/directional_derivative_mountain

Level Sets: Topographic Map

The function is the altitude



Topographic map



3-D picture

Level Set: Exercise

Consider a strictly convex-quadratic function

$$f(x) = \frac{1}{2}(x - x^*)^\top H(x - x^*) = \frac{1}{2} \sum_i h_{ii}(x_i - x_i^*) + \frac{1}{2} \sum_{i \neq j} h_{ij}(x_i - x_i^*)(x_j - x_j^*)$$

with H a symmetric, positive, definite matrix ($H \succ 0$).

1. What is/are the optima of f ? What does H represent for the function?

2. Assume $n=2$, $H = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ plot the level sets of f

3. Same question with $H = \begin{bmatrix} 1 & 0 \\ 0 & 9 \end{bmatrix}$

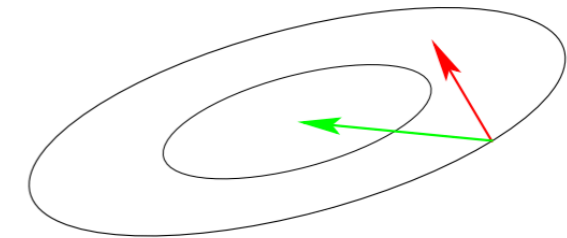
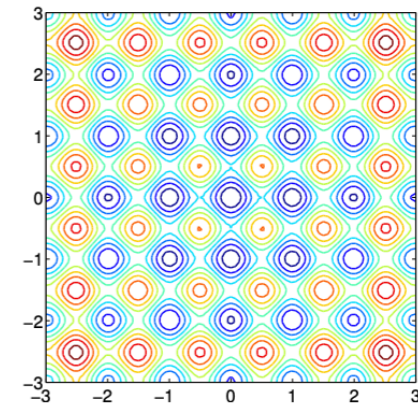
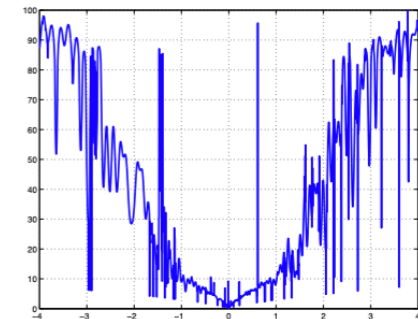
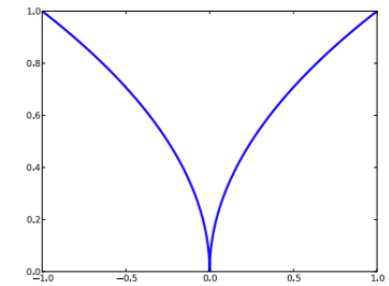
4. Same question with $H = P \begin{bmatrix} 1 & 0 \\ 0 & 9 \end{bmatrix} P^\top$ with $P \in \mathbb{R}^{2 \times 2}$
 P orthogonal

What Makes an Optimization Problem Difficult?

What Makes a Function Difficult to Solve?

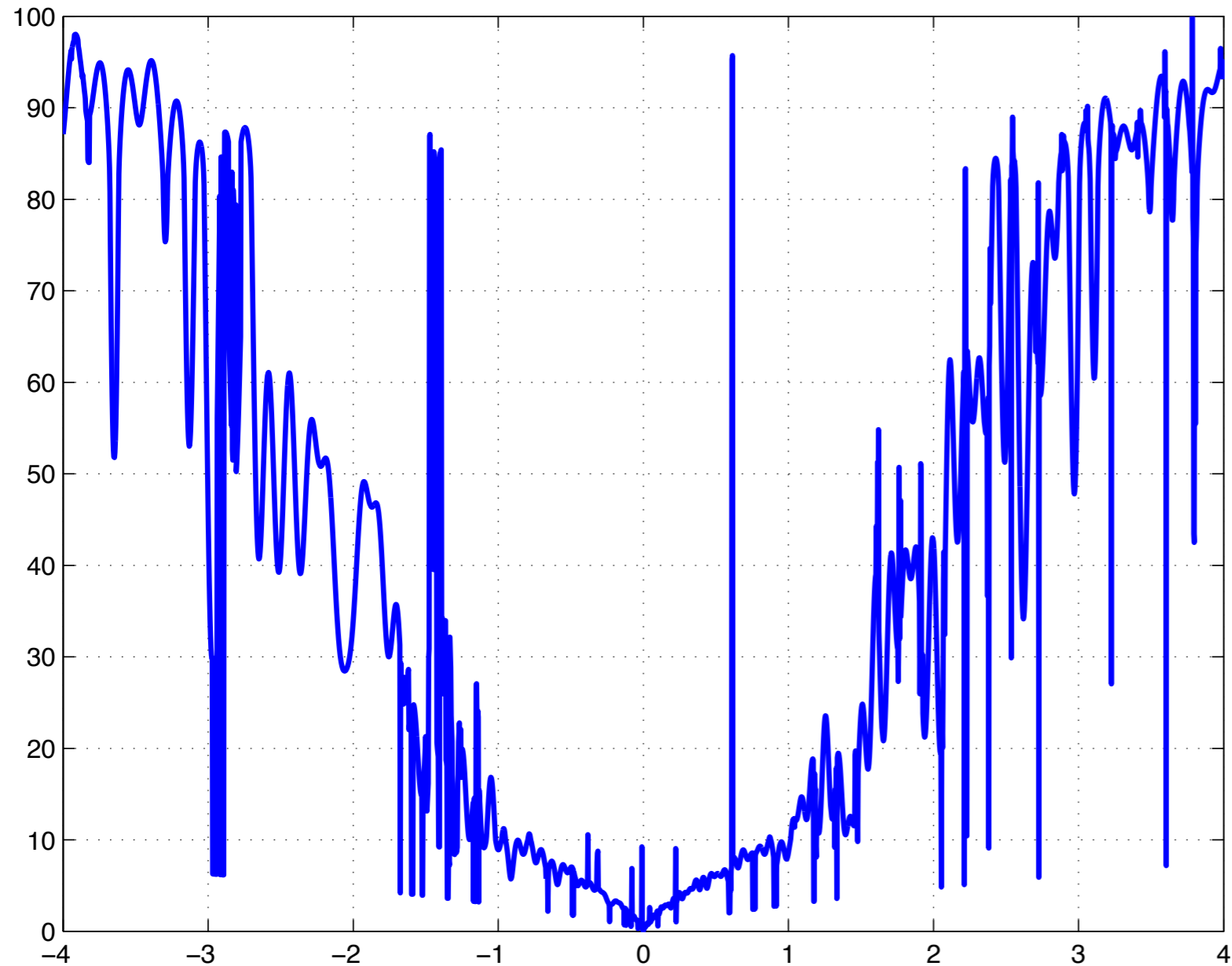
- ▶ non-linear, non-quadratic, non-convex
 - on linear and quadratic functions much better search policies are available
- ▶ ruggedness
 - non-smooth, discontinuous, multimodal, and/or noisy function
- ▶ dimensionality (size of search space)
 - (considerably) larger than three
- ▶ non-separability
 - dependencies between the objective variables
- ▶ ill-conditioning

Why stochastic search?



gradient direction Newton direction

Ruggedness



A cut of a 4-D function that can easily be solved with the
CMA-ES algorithm

Why is Optimization a non-trivial Problem?

Curse of dimensionality

if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

Why is Optimization a non-trivial Problem?

Curse of dimensionality

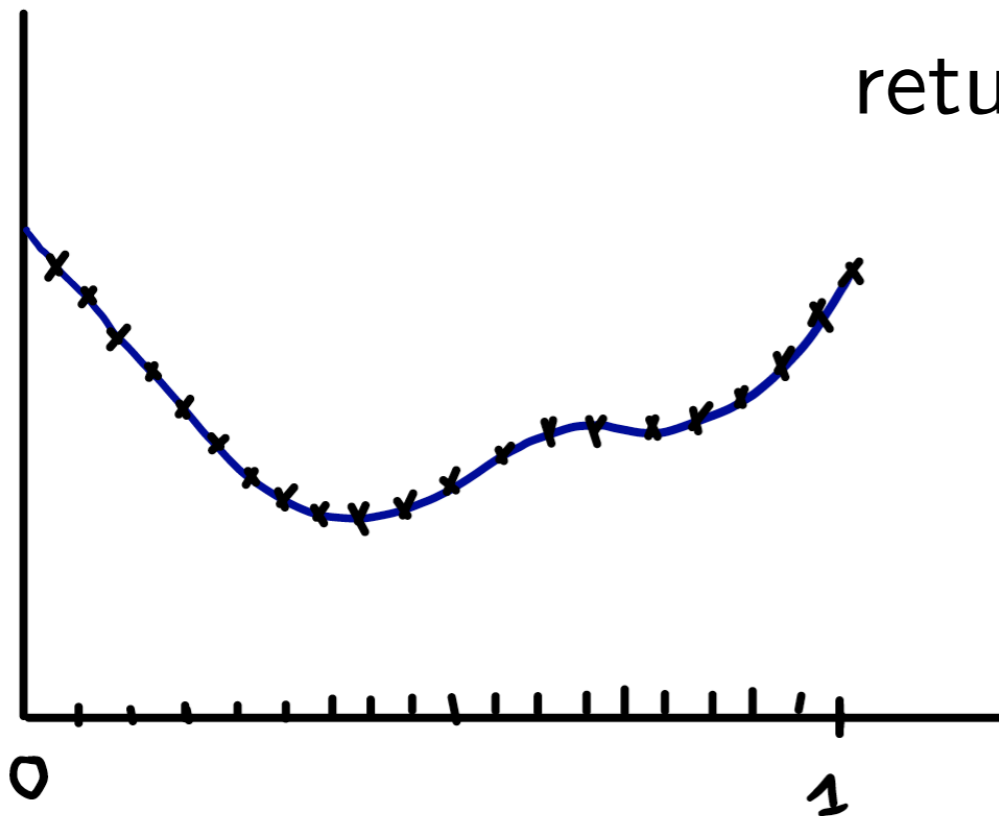
if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

set a regular grid on $[0,1]$

evaluate on f all the points of the grid

return the lowest function value



Why is Optimization a non-trivial Problem?

Curse of dimensionality

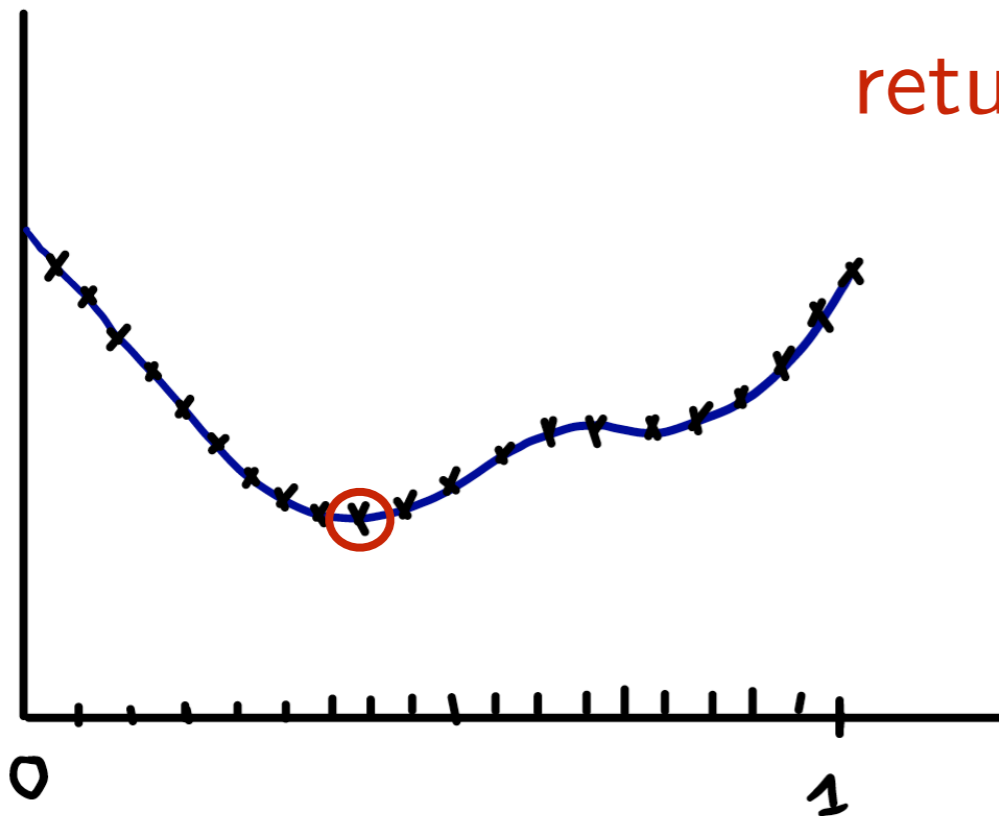
if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

set a regular grid on $[0,1]$

evaluate on f all the points of the grid

return the lowest function value



Why is Optimization a non-trivial Problem?

Curse of dimensionality

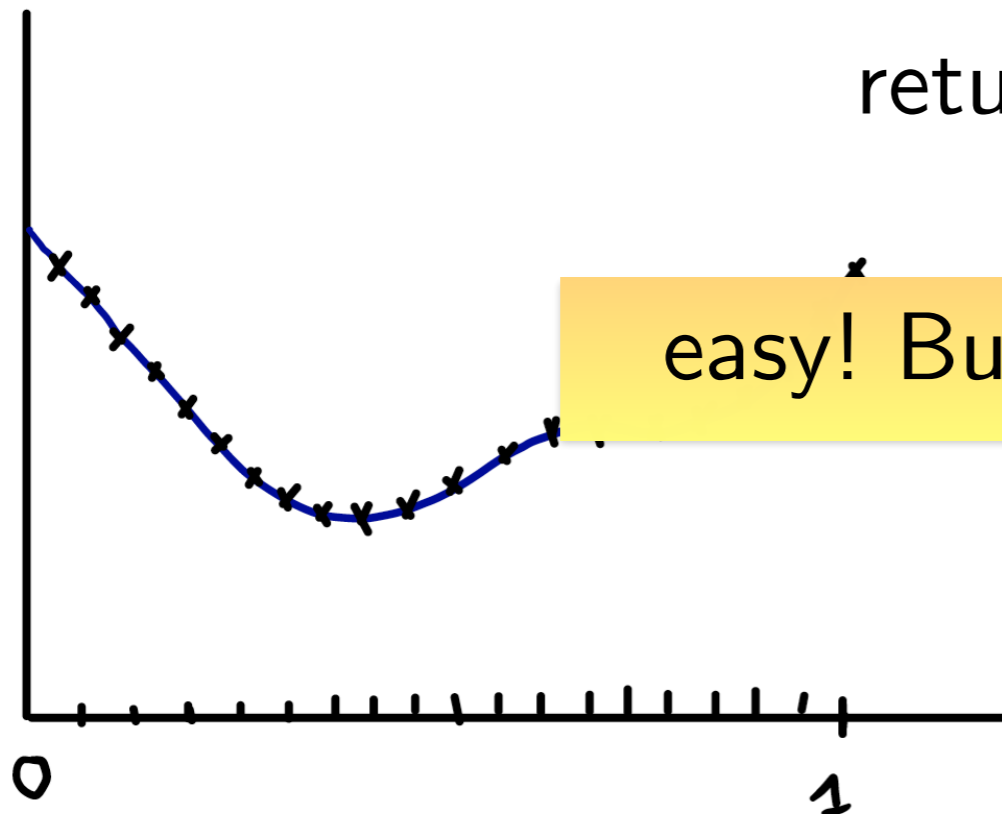
if $n=1$, which simple approach could you use to minimize:

$$f : [0, 1] \rightarrow \mathbb{R} \quad ?$$

set a regular grid on $[0,1]$

evaluate on f all the points of the grid

return the lowest function value



easy! But how does it scale when n increases?

1-D optimization is trivial

Curse of Dimensionality

The term **curse of dimensionality** (Richard Bellman) refers to problems caused by the **rapid increase in volume** associated with adding extra dimensions to a (mathematical) space.

Example: Consider placing 100 points onto a real interval, say $[0,1]$.

How many points would you need to get a similar coverage (in terms of distance between adjacent points) in dimension 10?

Curse of Dimensionality

The term **curse of dimensionality** (Richard Bellman) refers to problems caused by the **rapid increase in volume** associated with adding extra dimensions to a (mathematical) space.

Example: Consider placing 100 points onto a real interval, say $[0,1]$. To get similar coverage, in terms of distance between adjacent points, of the 10-dimensional space $[0,1]^{10}$ would require $100^{10} = 10^{20}$ points. A 100 points appear now as isolated points in a vast empty space.

Consequence: a **search policy** (e.g. exhaustive search) that is valuable in small dimensions **might be useless** in moderate or large dimensional search spaces.

Curse of Dimensionality

How long would it take to evaluate 10^{20} points?

Curse of Dimensionality

How long would it take to evaluate 10^{20} points?

```
import timeit
timeit.timeit('import numpy as np ;
np.sum(np.ones(10)*np.ones(10))', number=1000000)
> 7.0521080493927
```

7 seconds for 10^6 evaluations of $f(x) = \sum_{i=1}^{10} x_i^2$

We would need more than 10^8 days for evaluating 10^{20} points

[As a reference: origin of human species: roughly 6×10^8 days]

Separability

Given $x = (x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ denote

$$x^{-i} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \in \mathbb{R}^{n-1}$$

$$f_{x^{-i}}(y) = f(x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_n)$$

The function $f_{x^{-i}}(y)$ is a 1-D function which is a cut of f along the coordinate i .

Definition: A function f is **separable** if for all i , for all x, \bar{x}

$$\operatorname{argmin}_y f_{x^{-i}}(y) = \operatorname{argmin}_y f_{\bar{x}^{-i}}(y)$$

→ the optimum along the coordinate i , does not depend on how the other coordinates are fixed.

a weak definition of separability

Lemma: Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \text{Im}(f) \rightarrow \mathbb{R}$ strictly increasing. If f is **separable** then $g \circ f$ is separable.

Proposition: Let f be a **separable** then for all x

$$\operatorname{argmin} f(x_1, \dots, x_n) = \left(\operatorname{argmin}_y f_{x_{\neg 1}}(y), \dots, \operatorname{argmin}_y f_{x_{\neg n}}^n(y) \right)$$

and f can be optimized using n minimization along the coordinates.

Exercise: prove the proposition

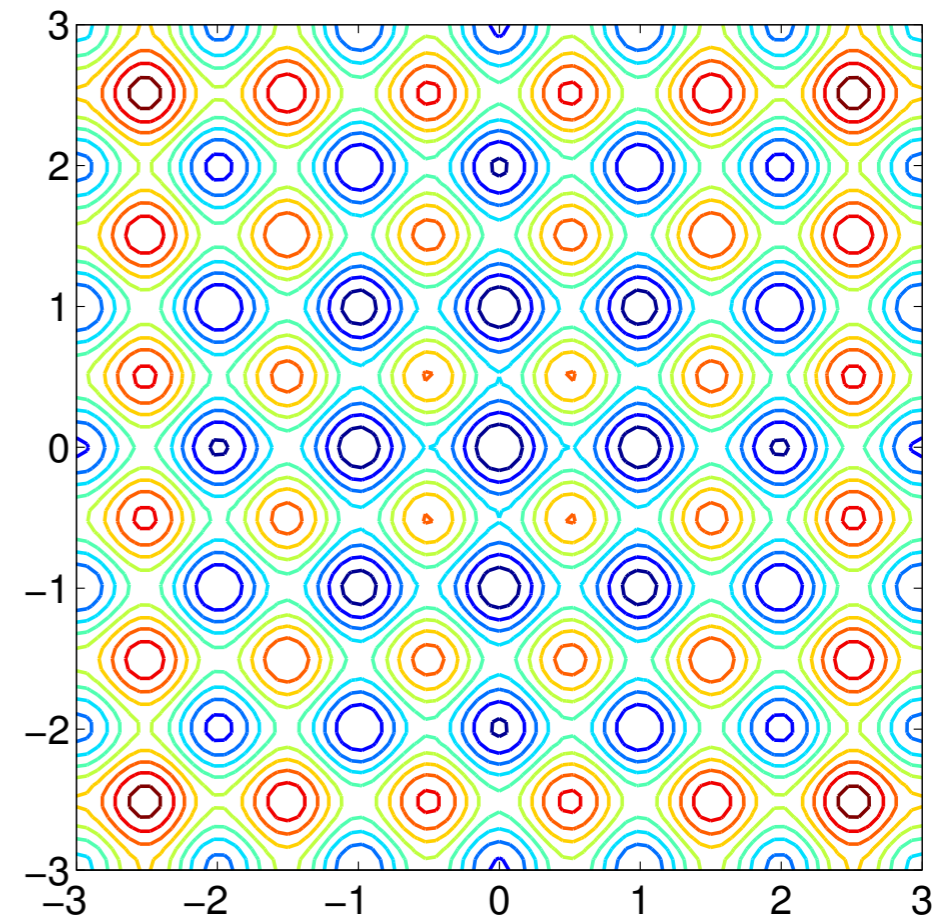
Example: Additively Decomposable Functions

Lemma: Let $f(x_1, \dots, x_n) = \sum_{i=1}^n h_i(x_i)$ for h_i having a unique argmin.

Then f is separable. We say in this case that f is additively decomposable.

Example: Rastrigin function

$$f(x) = 10n + \sum_{i=1}^n (x_i^2 - 10 \cos(2\pi x_i))$$



Consequence

Consider $f(x) = \prod_{i=1}^n h_i(x_i)$ with $h_i(x_i) > 0$. Then it is separable.

Non-separable Problems

Separable problems are typically easy to optimize. Yet **difficult real-world problems are non-separable.**

One needs to be careful when evaluating optimization algorithms that not too many test functions are separable and if so that the *algorithms do not exploit separability.*

***Otherwise:** good performance on test problems will not reflect good performance of the algorithm to solve difficult problems*

Algorithms known to exploit separability:

Many Genetic Algorithms (GA), Most Particle Swarm Optimization (PSO)

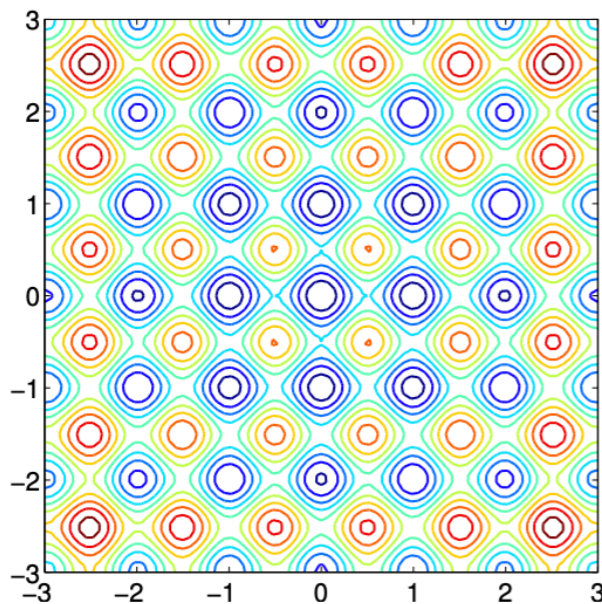
Non-separable Problems

Building a non-separable problem from a separable one

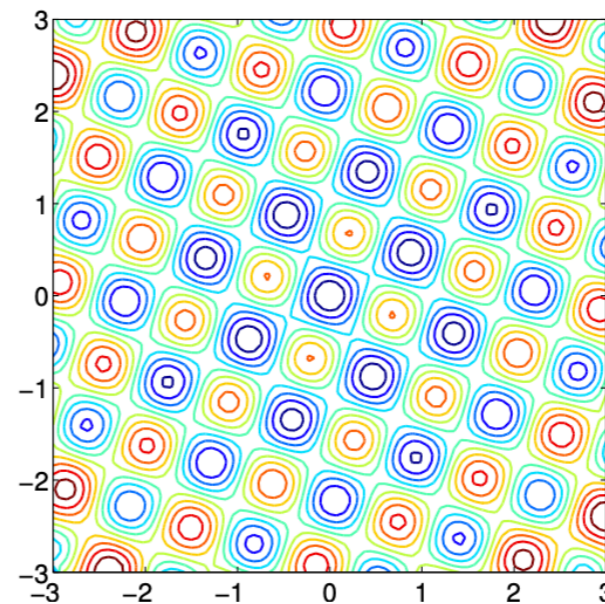
Rotating the coordinate system

- ▶ $f : \mathbf{x} \mapsto f(\mathbf{x})$ separable
- ▶ $f : \mathbf{x} \mapsto f(\mathbf{R}\mathbf{x})$ non-separable

R rotation matrix



R
→



¹ Hansen, Ostermeier, Gawelczyk (1995). On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation. Sixth ICGA, pp. 57-64, Morgan Kaufmann

² Salomon (1996). "Reevaluating Genetic Algorithm Performance under Coordinate Rotation of Benchmark Functions; A survey of some theoretical and practical aspects of genetic algorithms." BioSystems, 39(3):263-278

Ill-conditioned Problems - Case of Convex-quadratic functions

Consider a strictly convex-quadratic function

$$f(x) = \frac{1}{2}(x - x^*)^\top H(x - x^*) \text{ for } x = (x_1, \dots, x_n)^\top \in \mathbb{R}^n \text{ and}$$

$x^* \in \mathbb{R}^n$ with H a symmetric, positive, definite (SPD) matrix.

Remember that $H = \nabla^2 f(x)$.

The condition number of the matrix H (with respect to the Euclidean norm) is defined as

$$\text{cond}(H) = \frac{\lambda_{\max}(H)}{\lambda_{\min}(H)}$$

with $\lambda_{\max}()$ and $\lambda_{\min}()$ being respectively the largest and smallest eigenvalues.

Ill-conditioned means a high condition number of the Hessian matrix H .

Consider now the specific case of the function $f(x) = \frac{1}{2}(x_1^2 + 9x_2^2)$

1. Compute its Hessian matrix, its condition number
2. Plots the level sets of f , relate the condition number to the axis ratio of the level sets of f
3. Generalize to a general convex-quadratic function

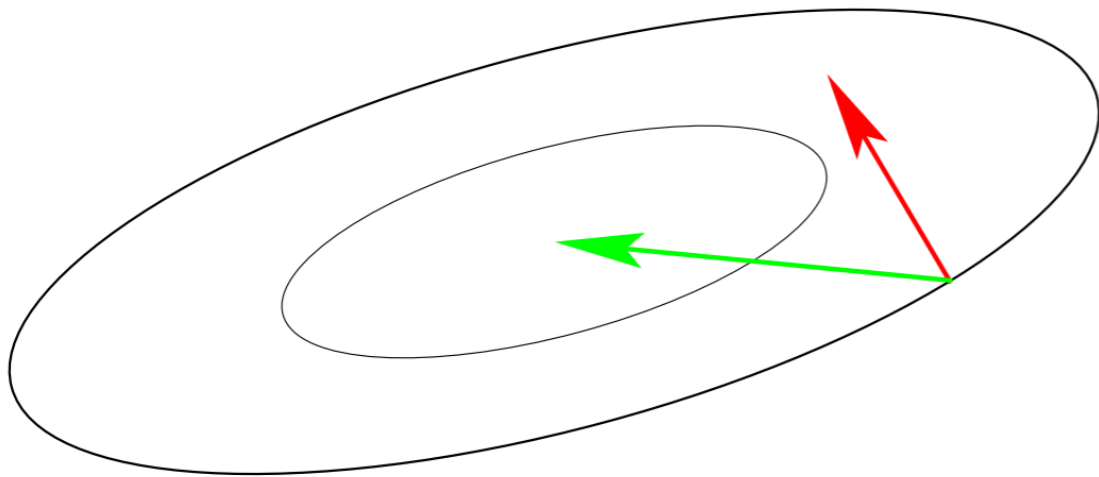
Real-world problems are often ill-conditioned.

4. Why do you think it is the case?
5. why are ill-conditioned problems difficult?

Ill-conditioned Problems

consider the curvature of the level sets of a function

ill-conditioned means “squeezed” lines of equal function value (high curvatures)



gradient direction $-f'(\mathbf{x})^T$

Newton direction
 $-H^{-1}f'(\mathbf{x})^T$

Condition number equals nine here. Condition numbers up to 10^{10} are not unusual in real world problems.