# Derivative Free Optimization

**joint course between
Optimization Master Paris Saclay - AMS Master
2025/2026**

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

| | |
|---|---|
| 28/11/2025 | room 1322 |
| 05/12/2025 | room 1322 |
| 12/12/2025 | room 1322 |
| 19/12/2025 | room 1322 |
| 09/01/2026 | room 1322 |
| 16/01/2026 | room 1322 |
| 23/01/2026 | room 1322 |
| 30/01/2026 | room 1322 |
| 06/02/2026 | room 1322 |
| **13/02/2026 [EXAM]** | TBA |

# Evaluation

**Written exam** on 14/02/2025

**Project (in group)** to be discussed

# Syllabus

**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

Benchmarking

**Practical Exercices**

practical exercices: 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

# Objectives

- understand how numerical optimization works when derivatives are unavailable
- identify the main challenges, typical failure modes
- learn how they are addressed in state-of-the-art methods
- teach about scientific / research approach in optimization

# Disclaimer & Focus

- This course does not aim at exhaustive coverage
- It emphasizes a major class of derivative-free methods: Evolution Strategies, with a strong focus on CMA-ES

# CMA-ES is:

one widely used derivative-free solvers worldwide
supported by Python implementations exceeding 80 million downloads
one of which will be used hands-on during the course
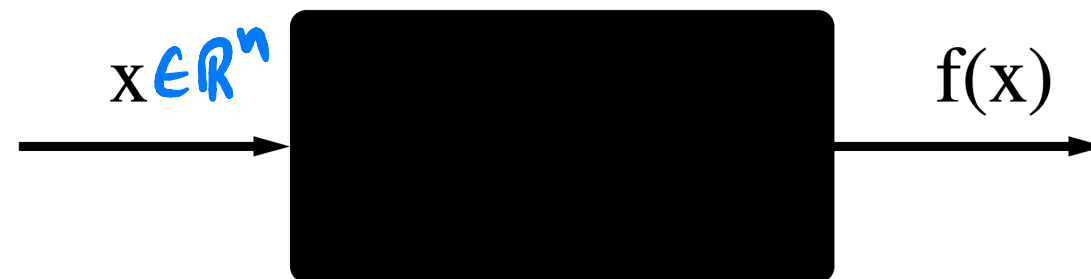
# Derivative-Free / Black-box Optimization

**Task:** minimize a numerical **objective** function (also called *fitness* function or *loss* function) $k \in \mathbb{N}$

$k = 1$ single-objective optimization

$k > 1$ : multi-objective optimization.

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

without derivatives (gradient). $\Omega$: search space, $n$ :dimension of the search space

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

$x \in \mathbb{R}^n$ → [black box] → f(x)

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

First order (black-box): methods that are using $Df$.

→ Gradient descent

→ Conjugate gradient

→ quasi-Newton (BFGS, ...)

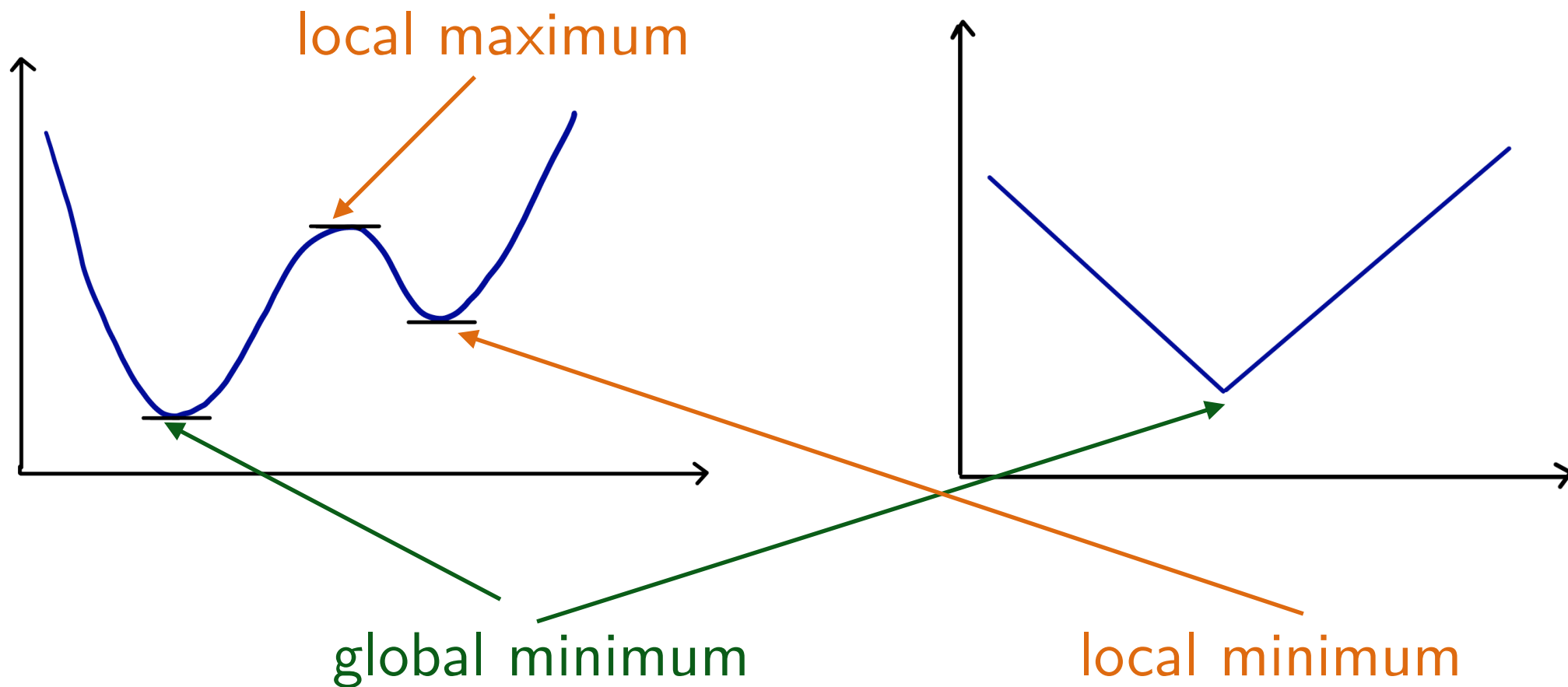Second order: $\nabla^2 f$, $Df$    Newton

<u>Newton direction:</u>    $-\nabla^2 f(x) \, Df(x)$

(followed by Newton-method)

n=1

local maximum

global minimum

local minimum

Poppy



Injection en orbite
- position
- vitesse

largage coiffe
(flux thermique)

flux thermique

120km

retombée d'étage

Séparations
(pression
dynamique)

fragmentation

visibilité

Vol atmosphérique
- efforts généraux
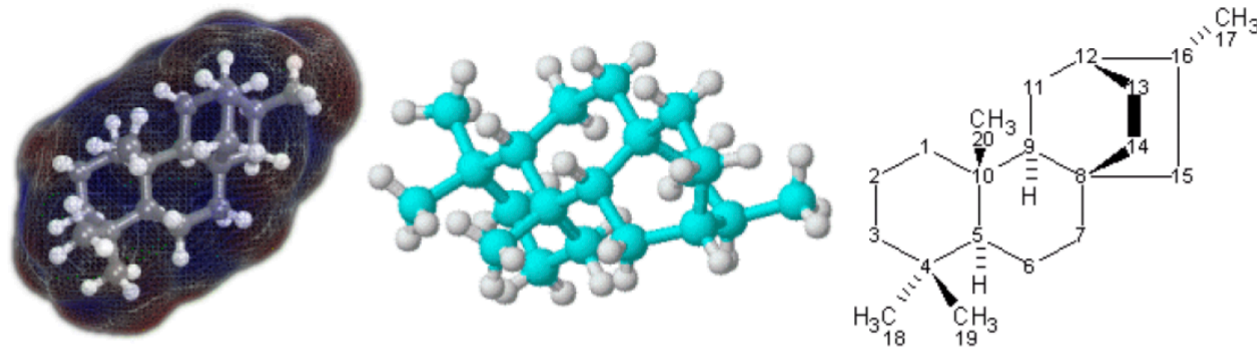- pilotage

pas de tir

station 1

station 2

copyright by Astrium

- 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
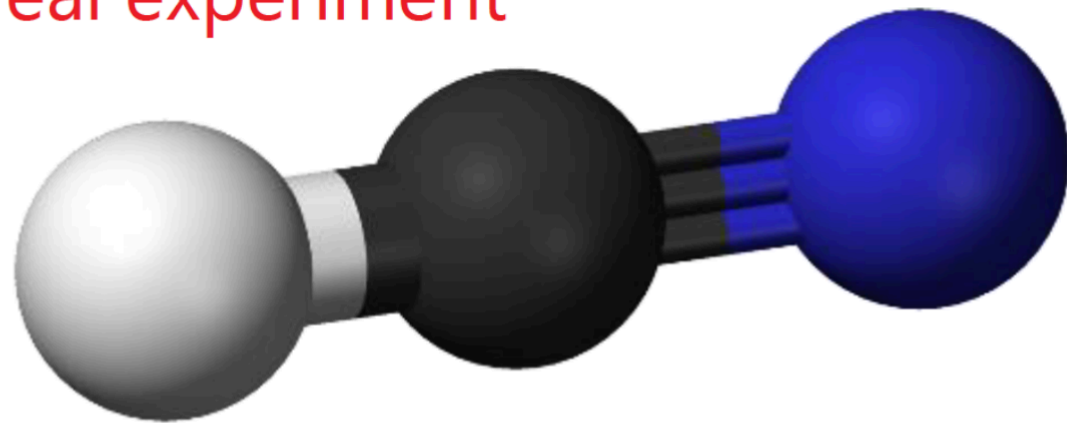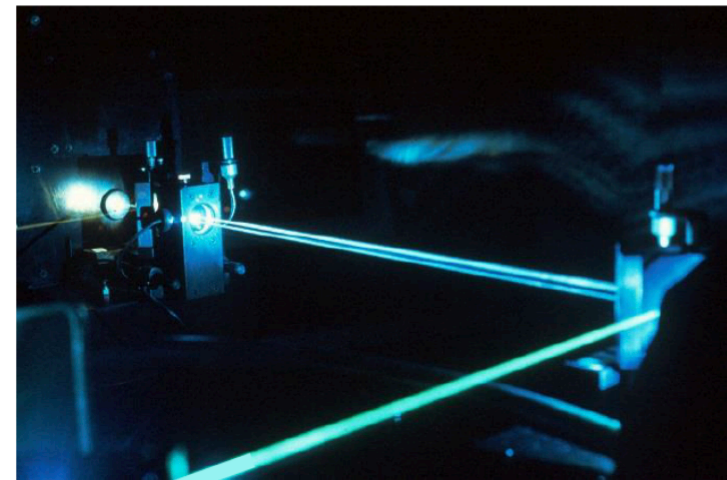
*23 continuous parameters to optimize*

*+ constraints*

# Control of the Alignement 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

$x_i \geq 0$    $x_i$ : quantity
$\sum x_i = 1$    of coffee $i$

▶ Find a mixture of coffee in order to keep the coffee taste from one year to another

▶ Objective function = opinion of one expert

$(x_1, \ldots, x_n) \longrightarrow (Taste(x)$  2
$\sim Taste\text{-}pair)$
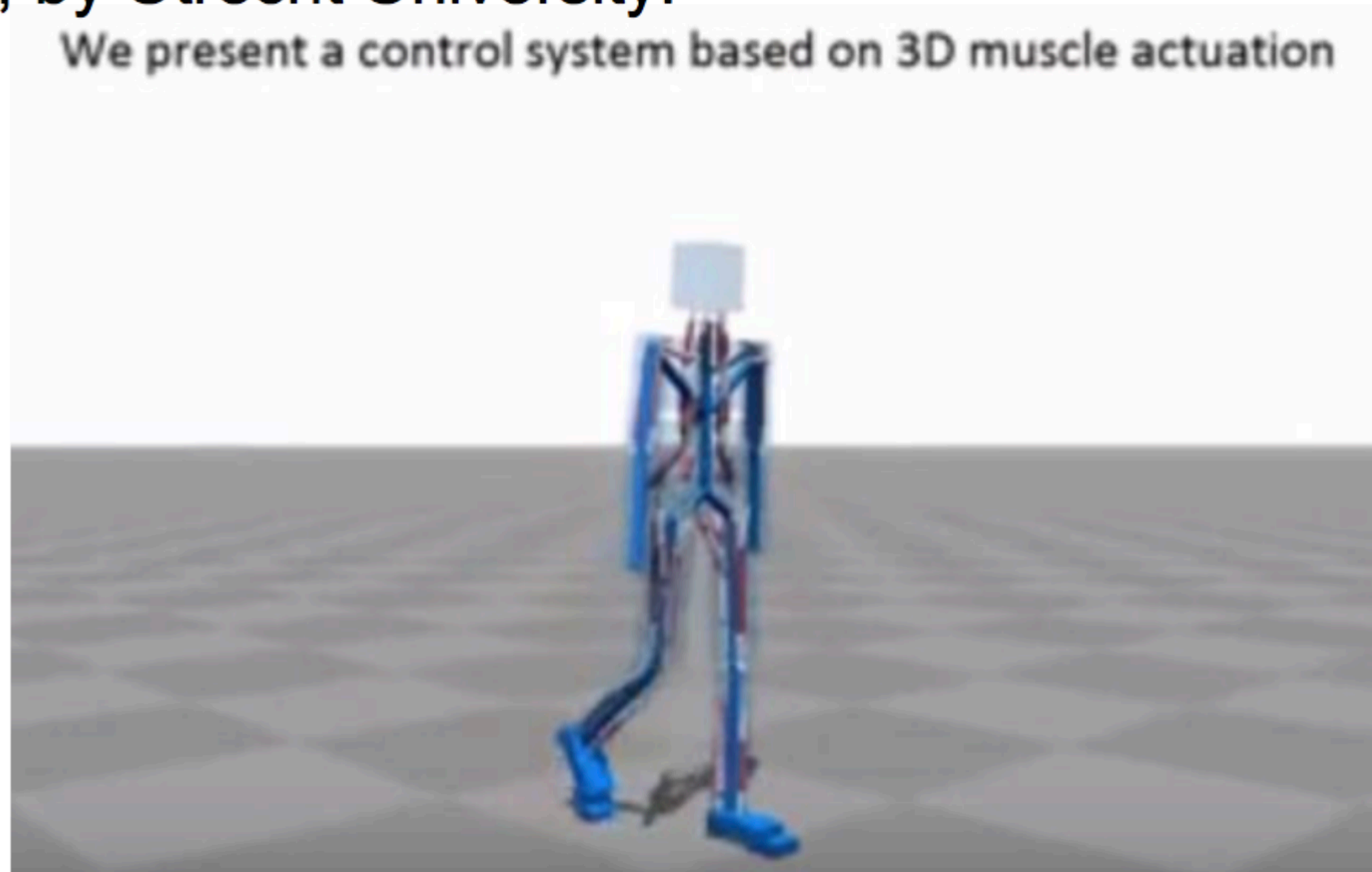$year)$



PUBLIC DOMAIN

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:



We present a control system based on 3D muscle actuation

https://www.youtube.com/watch?v=yci5Ful1ovk

T. Geitjtenbeek, 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 $x$

  $\hookrightarrow$ global minimizer

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

- 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 $\lim_{t \to \infty} f(x_t) = \min f$

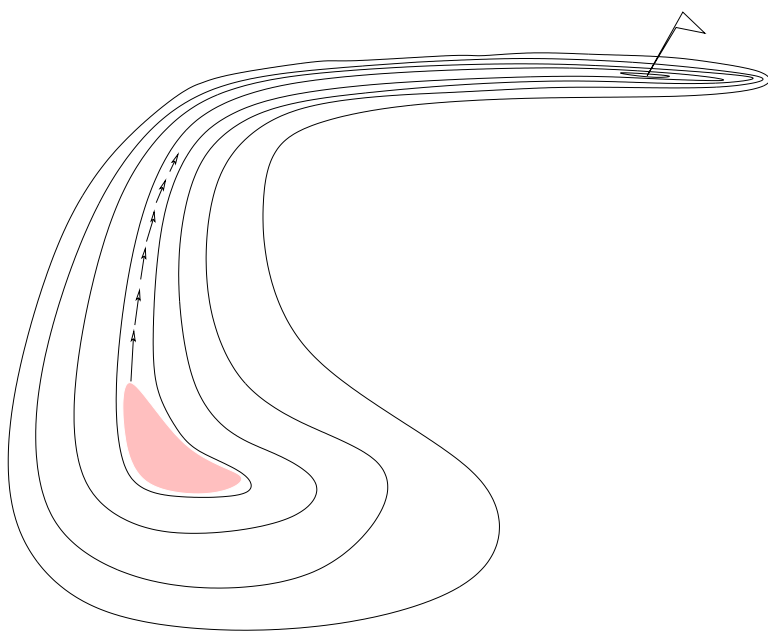  of course we want **fast** convergence

# Level Sets 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

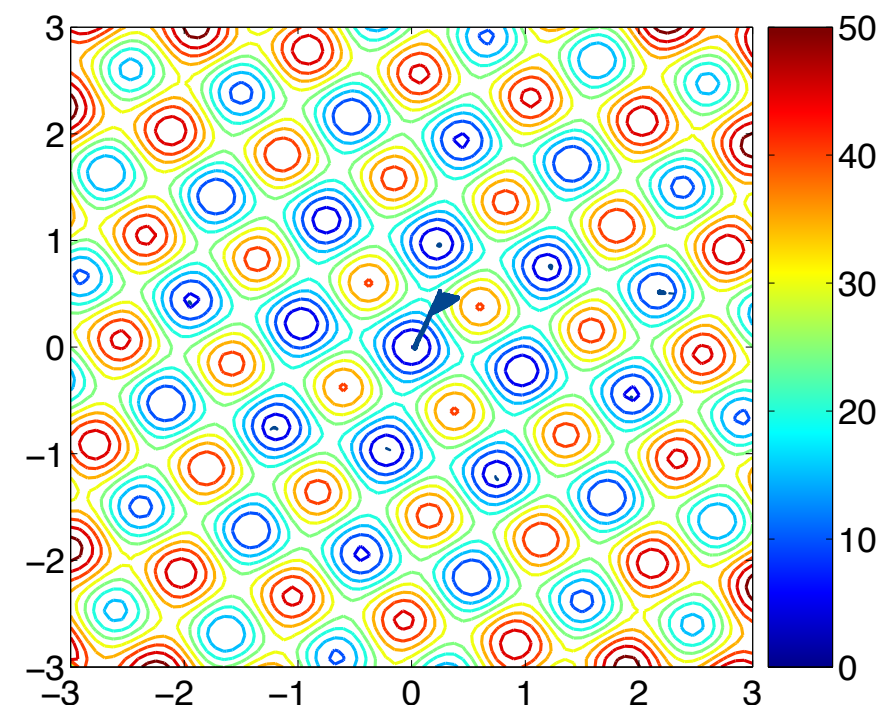$$\mathscr{L}_c = \{x \in \mathbb{R}^n \,|\, f(x) = c,\}, c \in \mathbb{R}$$

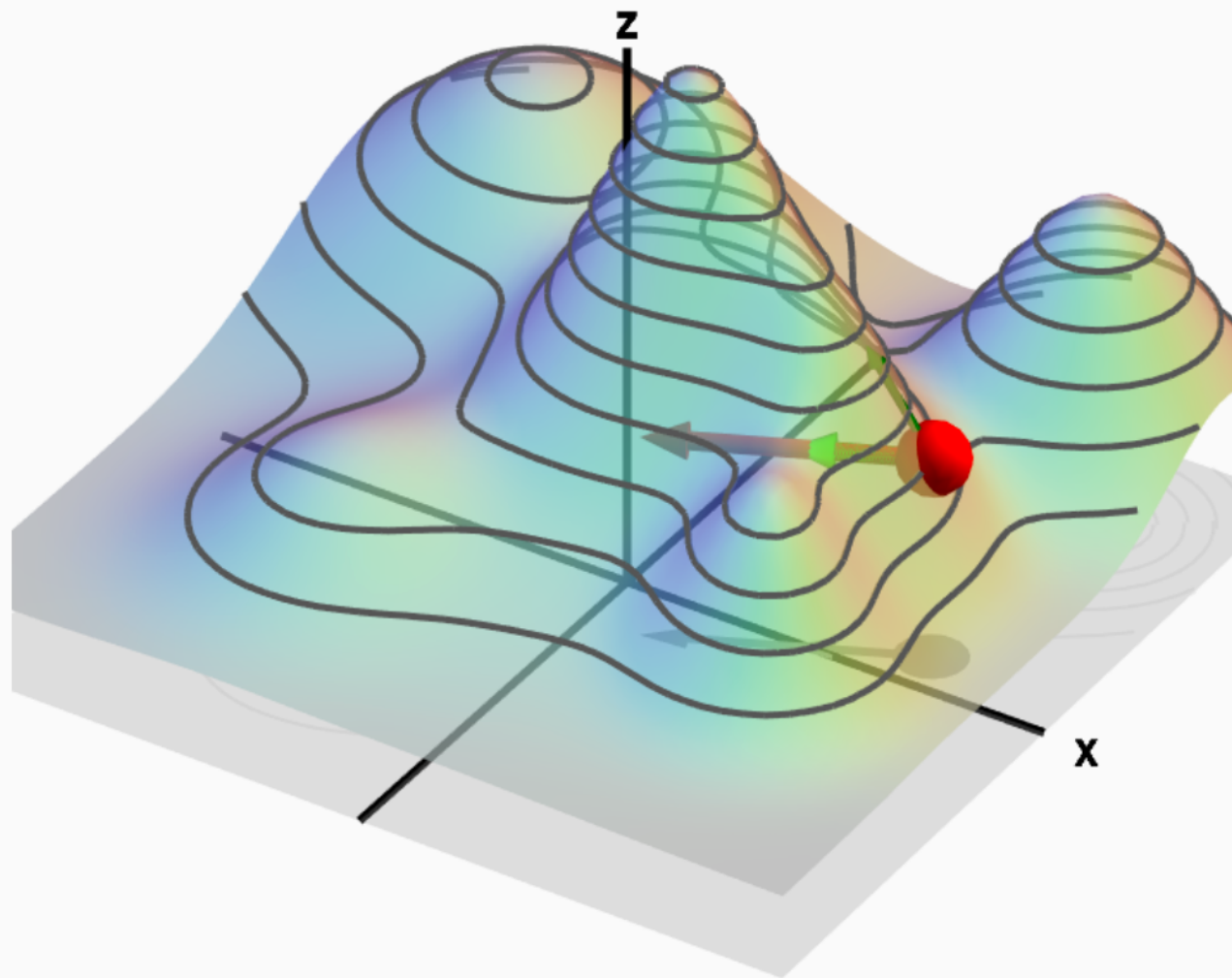Sublevel sets   $\text{Lev}_c^{\leq} = \{x \in \mathbb{R}^n \,|\, f(x) \leq c\}, \, c \in \mathbb{R}$

**Examples of level sets in 2D**

function with plenty local minima
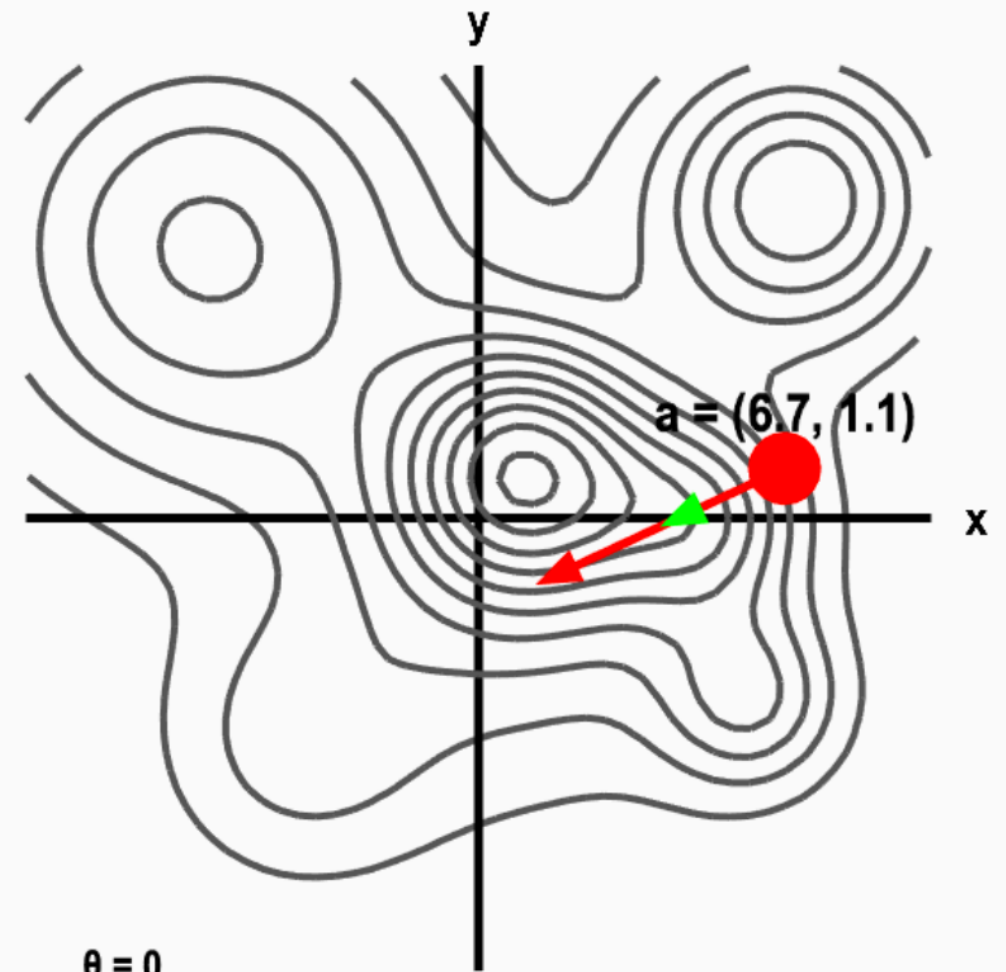
# Level Sets: Visualization of a Function



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

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

a = (6.7, 1.1)

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



3-D picture



Topographic map

Consider a strictly convex-quadratic function

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

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 $H = P \begin{bmatrix} 1 & 0 \\ 0 & 9 \end{bmatrix} P^T$ with $P \in \mathbb{R}^{2 \times 2}$

$P$ orthogonal

1/ $f(x) = \frac{1}{2}(x-x^a)^T H (x-x^a) \geq 0$   Since $H > 0$

$f(x^a) = 0$, thus $f(x^a) \leq f(x)$  $\forall x$

$\rightarrow x^a$ is the unique minimizer of $f$.

$f$ strictly convex.
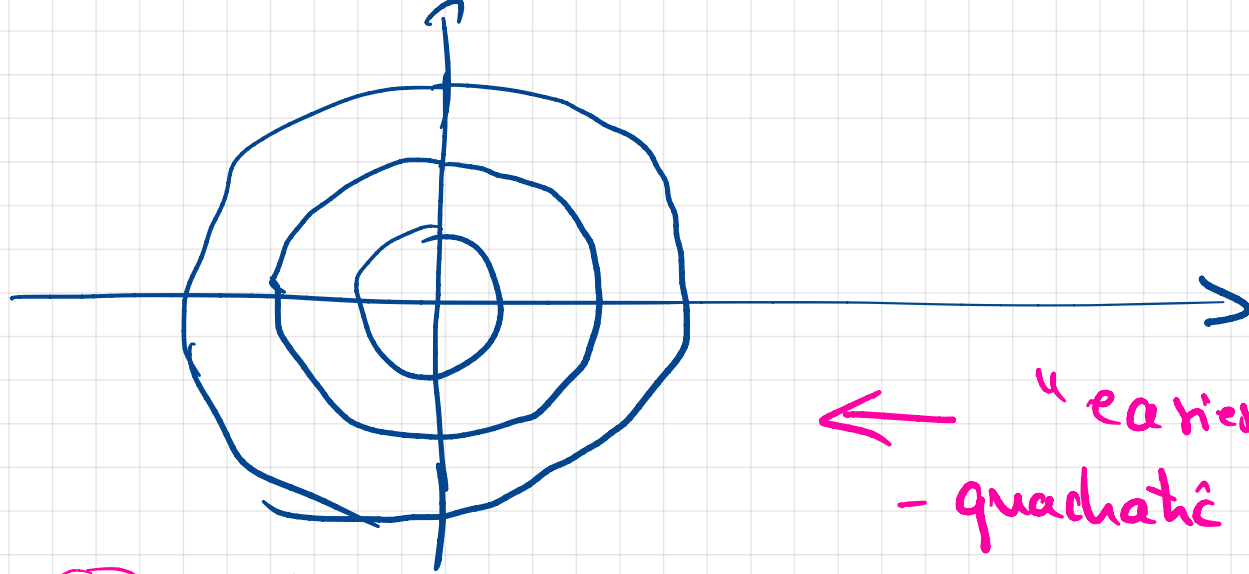
$H$ is the Hessian matrix. $D^2 f(x) = H$  $\forall x \in \mathbb{R}^n$

2/ $n = 2$,  $f(x) = \frac{1}{2}\left\{(x_1 - x_1^a)^2 + (x_2 - x_2^a)^2\right\}$

$\{x \mid f(x) = c\}$   $c > 0$

$= \{(x_1, x_2) \mid (x_1 - x_1^a)^2 + (x_2 - x_2^a)^2 = c\}$

"easiest" convex
– quadratic problem

3) If $H = \begin{pmatrix} \boxed{1} & \cdot \\ 0 & g \end{pmatrix}$

$$\text{lev}_{=c} f = \left\{ (x_1, x_2) \,\Big|\, \tfrac{1}{2}\left( (x_1 - x_1^*)^2 + g(x_2 - x_2^*)^2 \right) = c \right\}$$

$$c > 0$$

Assume WLG $x^* = 0$

$$\left\{ x = (x_1, x_2) \,\Big|\, \tfrac{1}{2}\left( x_1^2 + g\, x_2^2 \right) = c \right\}$$

$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$

$\begin{pmatrix} 3 \\ 0 \end{pmatrix}$

4) $H = P \begin{pmatrix} 1 & 0 \\ 0 & 9 \end{pmatrix} D^T \longrightarrow$ Rotated ellipsoid.

$P = \begin{pmatrix} p_1, & p_2 \end{pmatrix}$

$p_2$  $p_1$

Main- axis are the eigenvectors of $H$.

# What Makes an Optimization Problem Difficult?

# What makes a numerical optimization problem difficult?

1/ Non - convexity

2/ Non - smooth

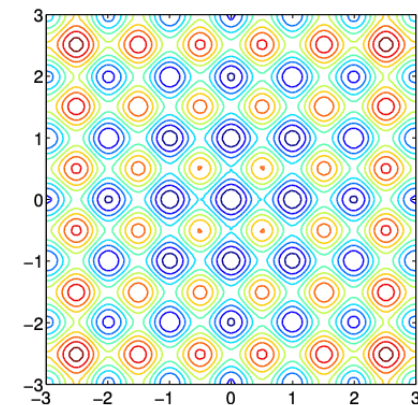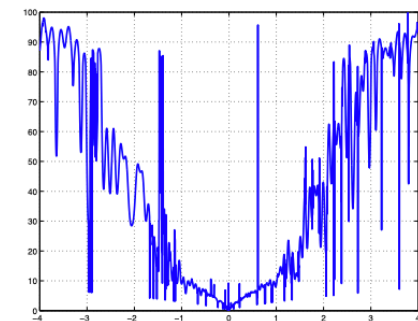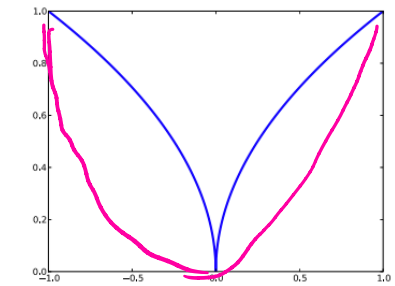3/ High - dimension     more than 100

4/ Steep : ill - conditionned

5/ Constraints

6/ "Expensive" to evaluate → Can sometimes take days (crash - simulations)
  ↳ parallelize
  ↳ "Bayesian optimization"

7/ If we do not know if there are optima
  ↳ It's typically the case : we do not ensure we found the global minimizer but we are happy with "good" local optimeum.

# What Makes a Function Difficult to Solve?

*Why stochastic search?*

▶ **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

↳ *more than one local optimum*

▶ **dimensionality (size of search space)**

(considerably) larger than three

*← many local optima*

▶ **non-separability**

dependencies between the objective variables

▶ **ill-conditioning**

**gradient direction** **Newton directio**

In this class:
strong focus on methods to address difficult black-box problems,
that can typically address the difficulties of the previous slide.

# Ruggedness



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

**Curse of dimensionality**

if n=1, which simple approach could you use to minimize:
$$f : [0, 1] \to \mathbb{R} \quad ?$$

## 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

**Curse of dimensionality**

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

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

set a regular grid on [0,1]

evaluate on f all the points of the grid

return the lowest function value

**Curse of dimensionality**
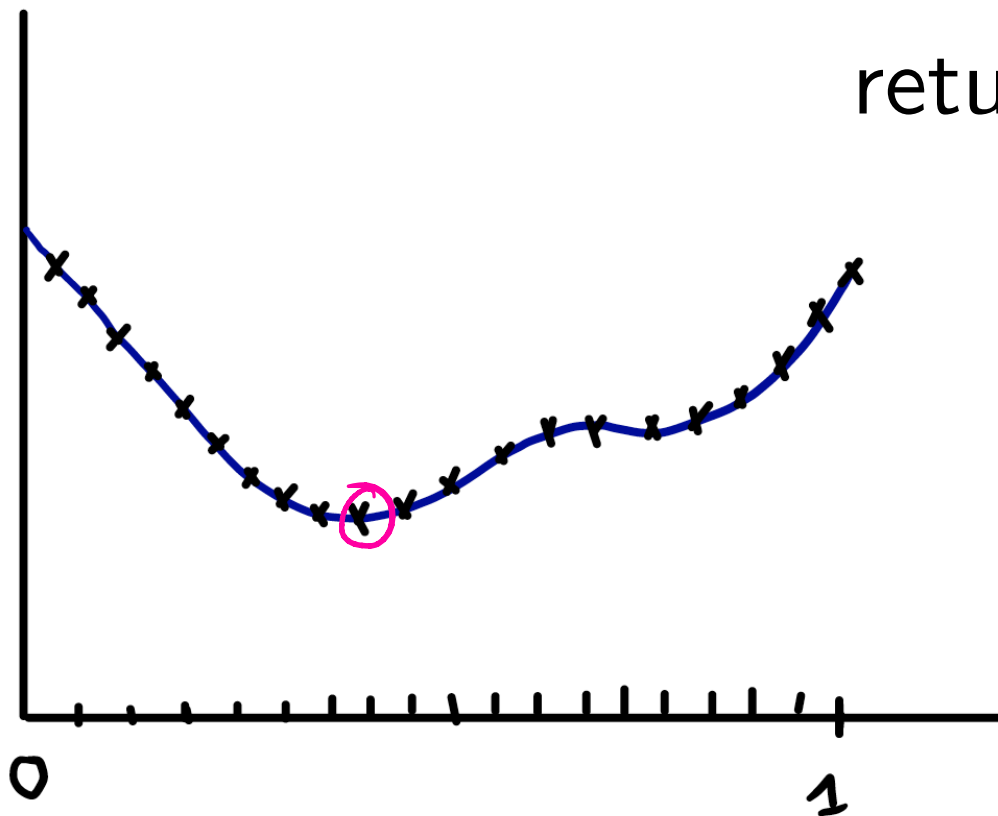
if n=1, which simple approach could you use to minimize:

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

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.

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

$$f(x) = \sum_{i=1}^{10} x_i^2$$

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
```

$$\sum_{i=1}^{n} x_i^2 \qquad x = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

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 x $10^8$ days]

# Separability

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

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

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

The function $f_{x^{\neg i}}(y)$ 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 $x, \bar{x}$

$$\mathrm{argmin}_y f_{x^{\neg i}}(y) = \mathrm{argmin}_y f_{\bar{x}^{\neg i}}(y)$$

$\rightarrow$ *the optimum along the coordinate i, does not depend on how the other coordinates are fixed.*

*a weak definition of separability*

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

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

$$\text{argmin} f(x_1, \ldots, x_n) = \left( \text{argmin}_y f_{x \neg 1}(y), \ldots, \text{argmin}_y f^n_{x \neg n}(y) \right)$$

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

**Exercice:** prove the proposition

# Example: Additively Decomposable Functions
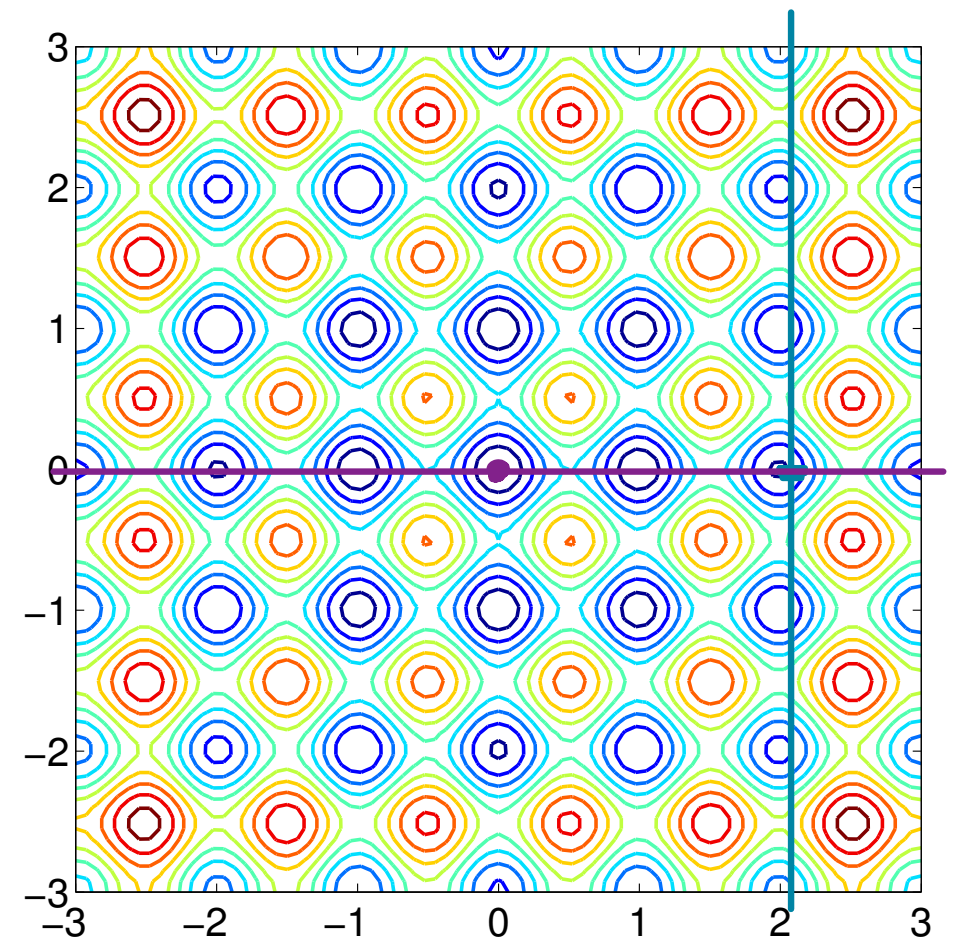
**Lemma:** Let $f(x_1, \ldots, x_n) = \sum_{i=1}^{n} h_i(x_i)$ for $h_i$ having a unique argmin.

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

# Non-separable Problems

Separable problems are typically easy to optimize. Yet **difficult real-word 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:

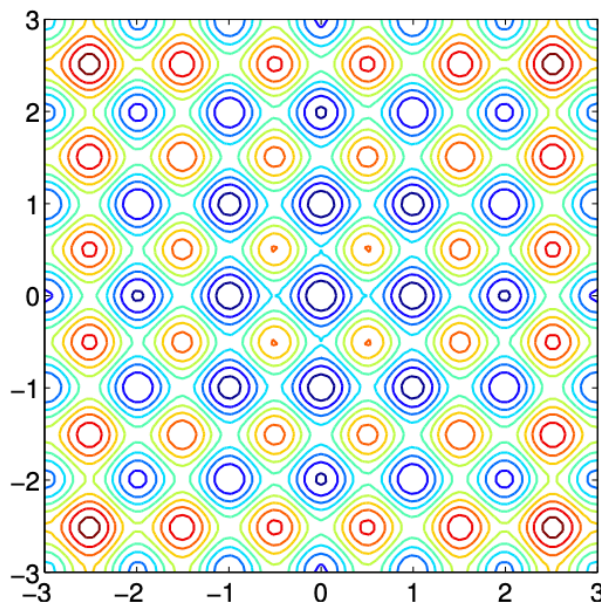Many Genetic Algorithms (GA), Most Particle Swarm Optimization (PSO)

# Non-separable Problems
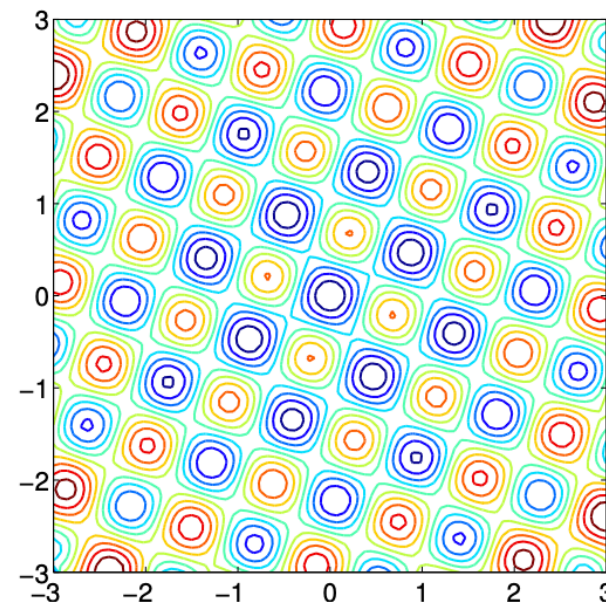
*Building a non-separable problem from a separable one*

## Rotating the coordinate system

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

$\mathbf{R}$ | *orthogonal matrix* / rotation matrix



$$\mathbf{R}$$

---

[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

Consider a strictly convex-quadratic function

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

$x^\star \in \mathbb{R}^n$ 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 $f(x) = \dfrac{1}{2}(x_1^2 + 9x_2^2)$

$H = \begin{pmatrix} 1 & 0 \\ 0 & 9 \end{pmatrix}$

**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 $f$

$cond(H) = 9$    $axis\text{-}ratio = 3 = \sqrt{cond(H)}$

**3.** Generalize to a general convex-quadratic function

Real-world problems are often ill-conditioned.

**4.** Why do you think it is the case?    → physical variables optimised can live on different scales.

**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 $-f'(\boldsymbol{x})^{\mathrm{T}}$

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

Ill-conditionned : cond $\sim 10^4/10^6$

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

# Part II: Algorithms

# A simple | randomized Derivative Free Optimization algorithm?
## | Stochastic

## PURE RANDOM SEARCH (PRS)

[Objective: minimize $f : [-1, 1]^n \to \mathbb{R}$

$X_t$ is the estimate of the optimum at iteration $t$

`Input` $(U_t)_{t \geq 1}$ independent identically distributed each $U_t \sim \mathcal{U}_{[-1,1]^n}$ (unif. distributed in $[-1, 1]^n$) ]

```
1. Initialize t = 1, X₁ = U₁
2. while not terminate
3.         t = t + 1
4.         If f(Uₜ) ≤ f(Xₜ₋₁)
5.                 Xₜ = Uₜ
6.         Else
8.                 Xₜ = Xₜ₋₁
```

1. Show that for all $t \geq 1$
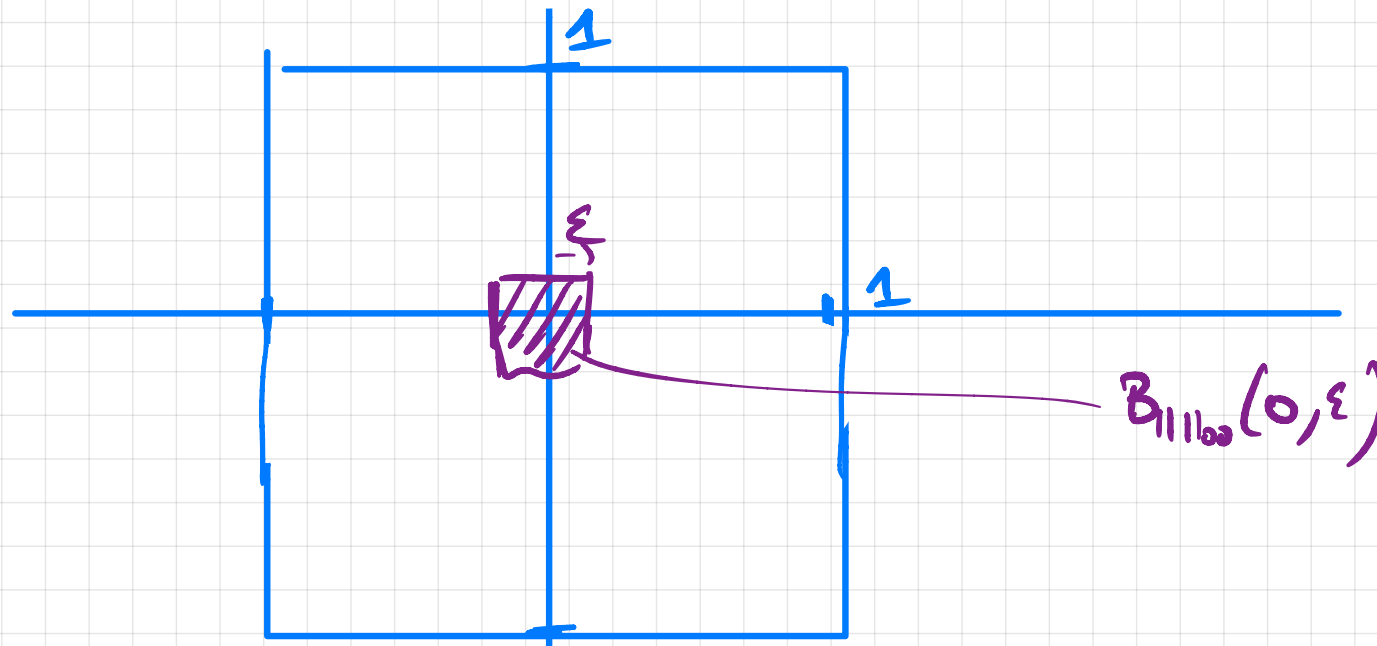
$$f(X_t) = \min\{f(U_1), \ldots, f(U_t)\}$$

## BY INDUCTION

2. We consider the simple case where $f(x) = \|x\|_\infty$ (we remind that $\|x\|_\infty := \max(|x_1|, \ldots, |x_n|)$). Show the convergence in probability of the PRS algorithm towards the optimum of $f$, that is prove that for all $\epsilon > 0$

$$\lim_{t \to \infty} \Pr\left(\|X_t\|_\infty \geq \epsilon\right) = 0$$

Hint: Prove and use the equality

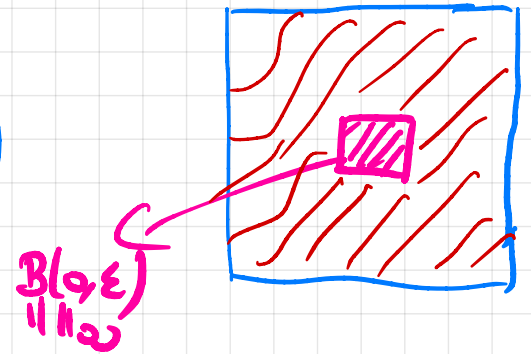$$\{\|X_t\|_\infty \geq \epsilon\} = \cap_{k=1}^{t}\{\|U_k\|_\infty \geq \epsilon\}$$



We want to prove that $\Pr\left(\|x_t\| \geq \epsilon\right) \xrightarrow[t \to +\infty]{} 0$

$$\left\{ \|X_t\|_\infty \geq \varepsilon \right\} \underset{(1)}{=} \left\{ \min \left\{ \|U_1\|_\infty, \ldots, \|U_t\|_\infty \right\} \geq \varepsilon \right\}$$

$$= \bigcap_{k=1}^{t} \left\{ \|U_k\|_\infty \geq \varepsilon \right\}$$

$$\Pr\left( \left\{ \|X_t\|_\infty \geq \varepsilon \right\} \right) = \Pr\left( \bigcap_{k=1}^{t} \left\{ \|U_k\|_\infty \geq \varepsilon \right\} \right)$$

$$\underset{\text{By ind}}{=} \prod_{k=1}^{t} \Pr\left( \|U_k\|_\infty \geq \varepsilon \right)$$

$$\underset{\substack{\text{Since the RV} \\ \text{are identically} \\ \text{dis}}}{=} \Pr\left( \|U_1\|_\infty \geq \varepsilon \right)^{t}$$

$$\Pr\left( \|U_1\|_\infty \geq \varepsilon \right) = 1 - \Pr\left( \|U_1\|_\infty \leq \varepsilon \right)$$

$= \Pr(\text{to be in red Set})$

$$= 1 - \frac{\text{Vol}\left( \left\{ x \mid \|x\|_\infty \leq \varepsilon \right\} \right)}{\text{Vol}\left( \left\{ x \mid \|x\|_\infty \leq 1 \right\} \right)}$$



$\frac{B(q,\varepsilon)}{\|\cdot\|_\infty}$

$$= 1 - \frac{(2\varepsilon)^n}{2^n}$$

$$= 1 - \varepsilon^n$$

$$\Pr\left(\left\{ \|X_t\|_\infty \geq \varepsilon \right\}\right) = \left(1 - \varepsilon^m\right)^t \xrightarrow[t \to \infty]{} 0$$

Hence the $\mathbb{P}RS$ converges in probability to the the minimizer of $f(x) = \|x\|_\infty$.

Let's quantify how fast it converges.

$\hookrightarrow$ for this we look at hitting time of the optimum.

$$T_\varepsilon = \inf \{ t , \|X_t\|_\infty \leq \varepsilon \}$$ Hitting time of

a ball $\wedge$ around the optimum.
   of radius $\varepsilon$.

$T_\varepsilon$ is a random variable

We can try to estimate the expected hitting time

$$\mathbb{E}(T_\varepsilon).$$

Let us assume we play a $\wedge$ game with 2 outcomes:
random

Win      with probability $p$

loose    _____  $1-p$

$$\left( \begin{array}{c} \text{Game} \\ \text{Head}/\text{Tail} \\ \rightarrow \text{proba to} \\ \text{win } \frac{1}{2} \end{array} \right)$$

and the outcome is sampled randomly and idependently.

We play till we win.

$T$ is how many times I need to play to see the first win.

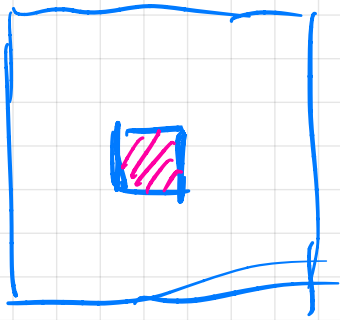$\rightarrow$ $T$ follows a geometric distribution with parameter $p$

On "average"      $\mathbb{E}(T) = \frac{1}{p}$

We can make a parallel between PRS and a game :

Win = if I reach $B(0, \varepsilon)$

lose otherwise.



$T_\varepsilon$ = Hitting time of $B(0, \varepsilon)$

= time it takes to win in the game.

$T_\varepsilon \sim$ geometric distribution with parameter

$$p = \Pr\left( \|U_1\|_\infty \leq \varepsilon \right) = \varepsilon^n$$

$$\mathbb{E}(T_\varepsilon) = \frac{1}{\varepsilon^n}$$

This is slow, as a comparison, if we set linear (geometric) convergence then

$$\mathbb{E}(T_\varepsilon) \approx \log(1/\varepsilon)$$

OBJECTIVE: Explain how to do better than PRS.