

Derivative Free Optimization

**joint course between
Optimization Master Paris Saclay - AMS
Master
2024/2025**

Anne Auger - Dimo Brockhoff
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

29/11/2024	room 1314
06/12/2024	room 1314
13/12/2024	room 1314
20/12/2024	room 1213
10/01/2025	room 1213
17/01/2025	room 1314
24/01/2025	room 1314
31/01/2025	room 1314
07/02/2025	room 1314
14/02/2025 [EXAM]	TBA

Evaluation

- Written exam on 14/02/2025) 60%
- Project (in group) around benchmarking of algorithms) 40%
 - oral presentation to the class

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 [taught D. Brockhoff]

Benchmarking (partly taught by D. Brockhoff)

Practical Exercises

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

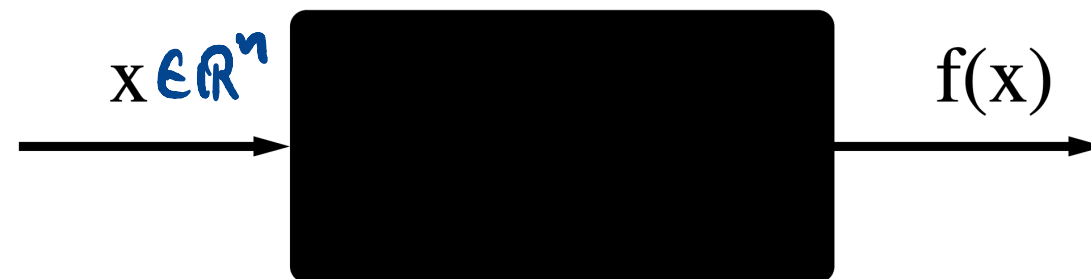
Derivative-Free / Black-box Optimization

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

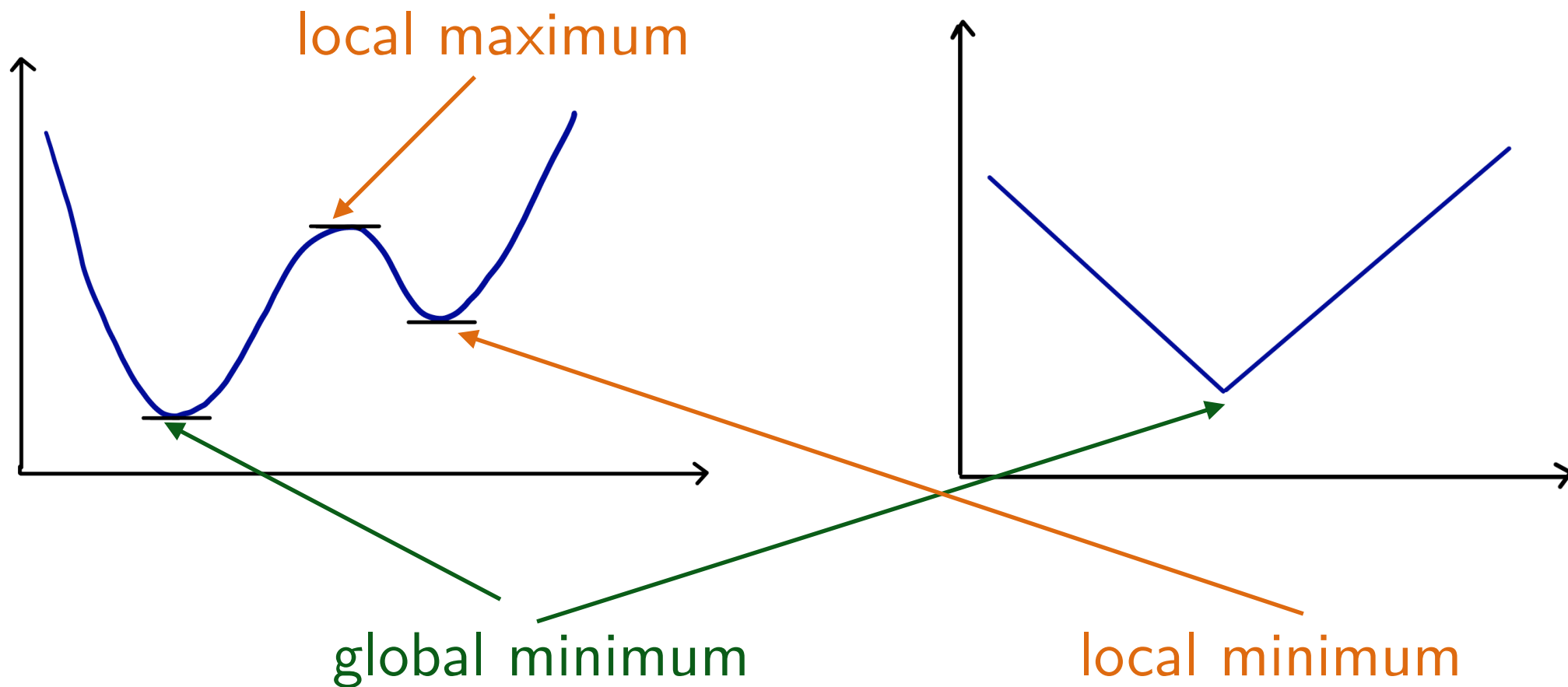
Also called **zero-order black-box** optimization



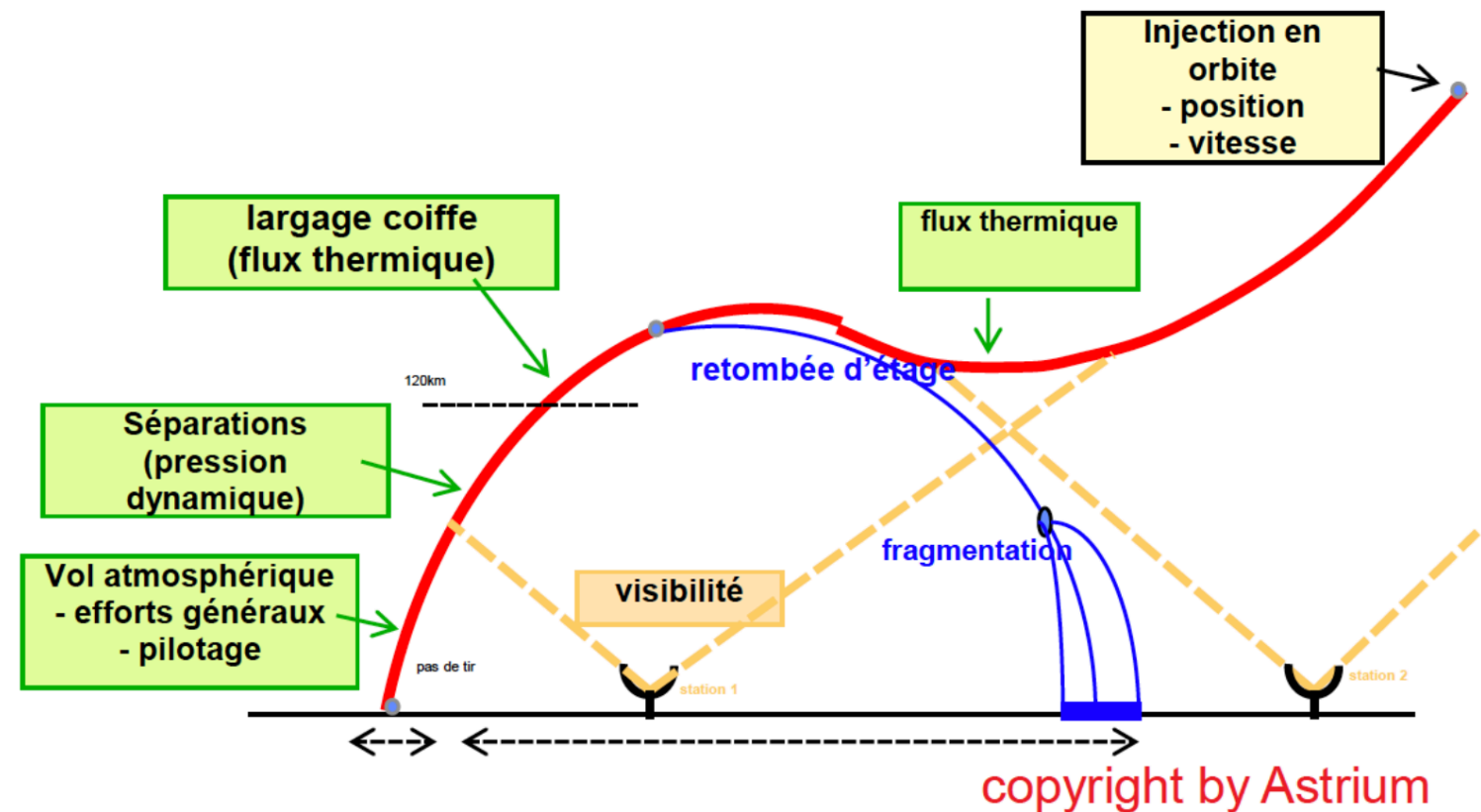
The function is seen by the algorithm as a zero-order **oracle** [a first order oracle would also return gradients] that can be queried at points and the oracle returns an answer

Reminder: Local versus Global Optimum

$n=1$



Examples: Optimization of the Design of a Launcher

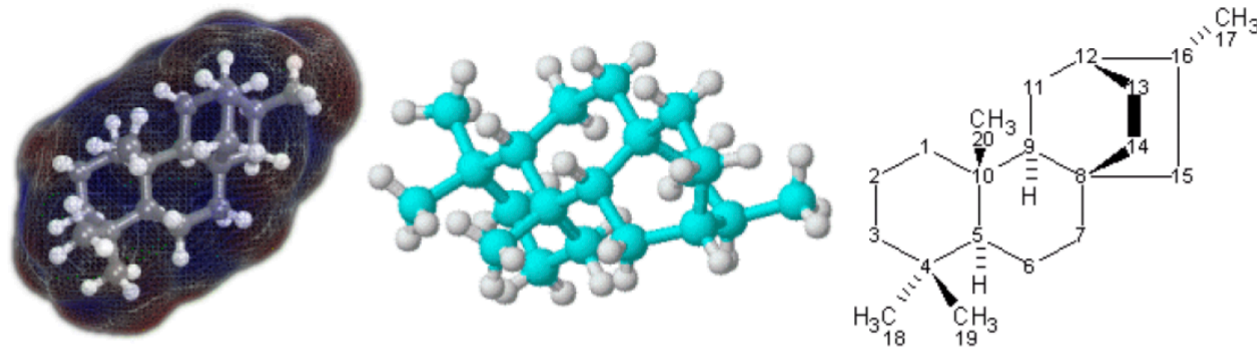


- Scenario: multi-stage launcher brings a satellite into orbit
- Minimize the overall cost of a launch
- Parameters: propellant mass of each stage / diameter of each stage / flux of each engine / parameters of the command law

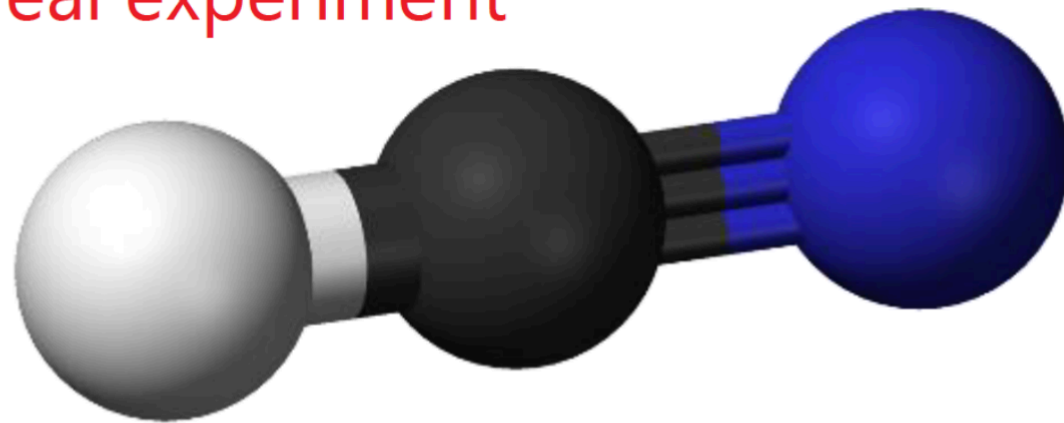
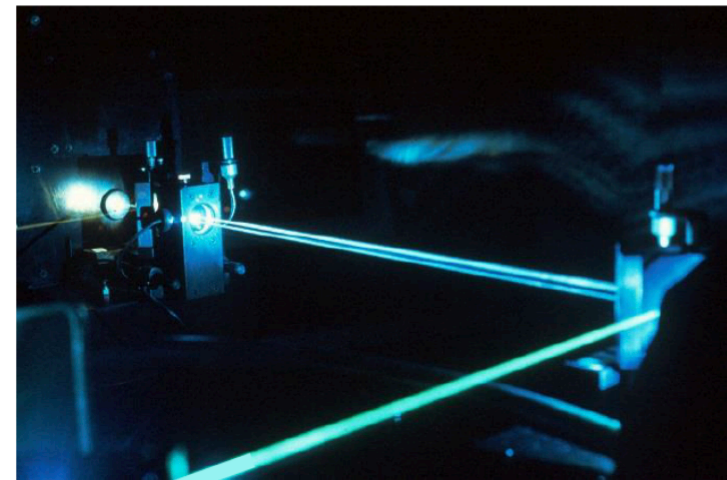
*23 continuous parameters to optimize
+ constraints*

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

$$x_i \geq 0$$
$$\sum x_i = 1$$

$$(x_1, x_2, x_3, x_4) \longrightarrow \text{Taste}$$

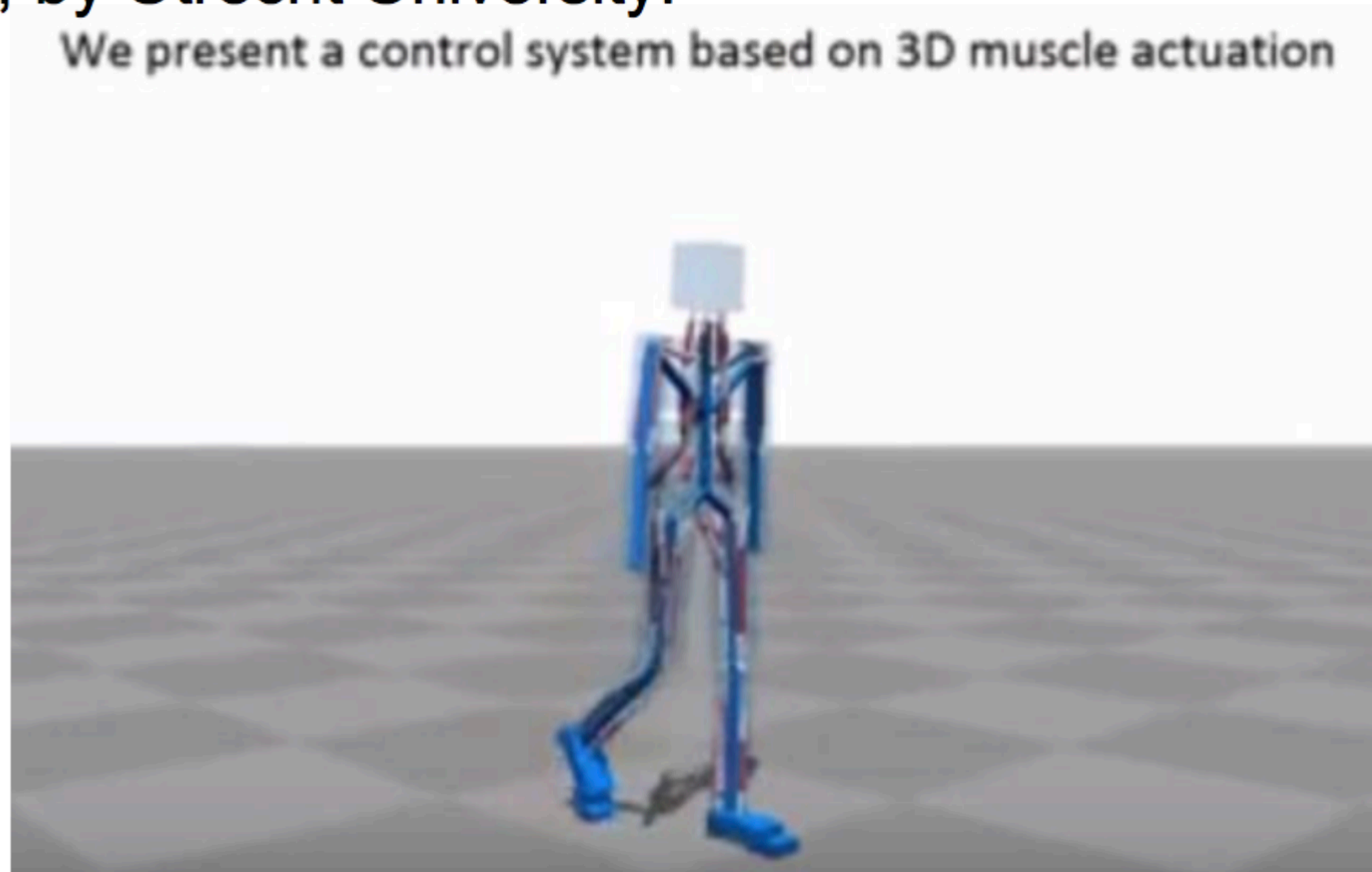


Quasipalm

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=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

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

- 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

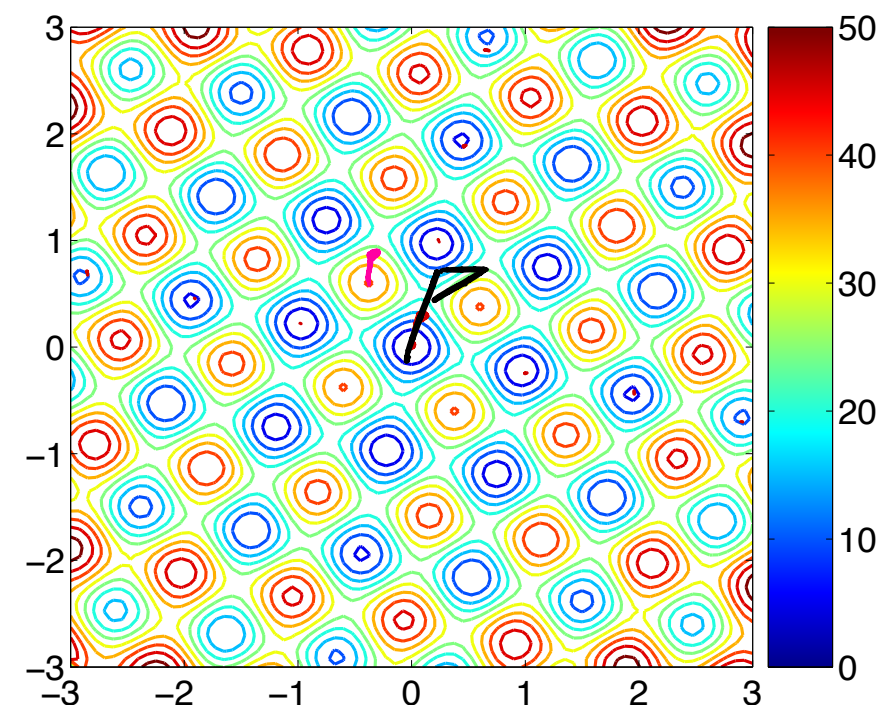
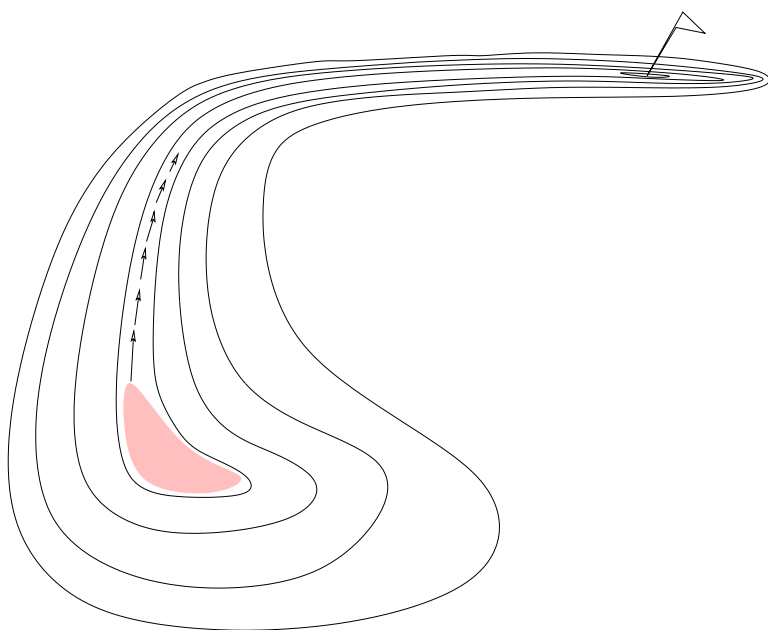
Level Sets of a Function

Level Sets: Visualization of a Function

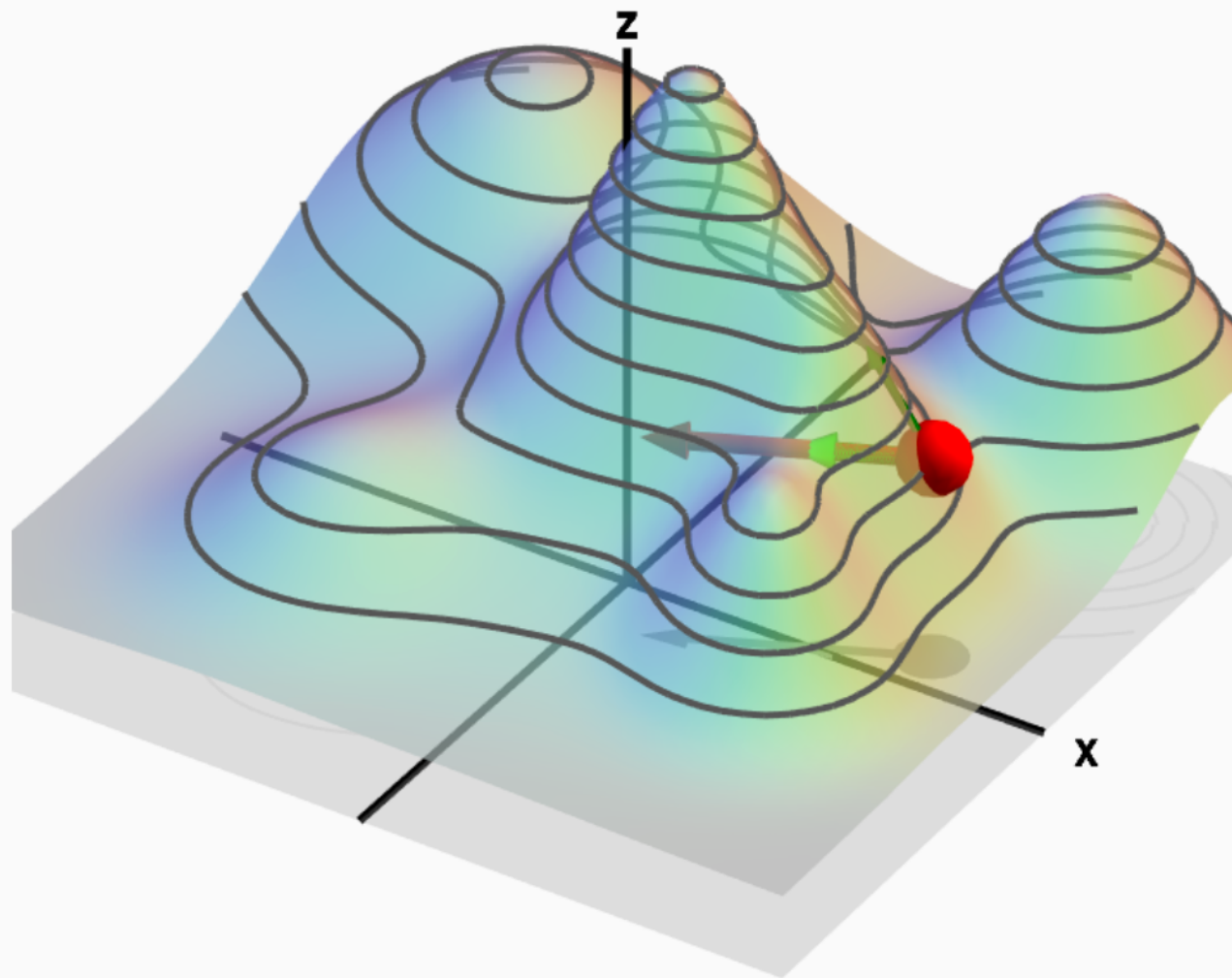
One-dimensional (1-D) representations are often misleading (as 1-D optimization is “trivial”, see slides related to curse of dimensionality), we therefore often represent **level-sets** of functions

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

Examples of level sets in 2D



Level Sets: Visualization of a Function



$\theta = 0$



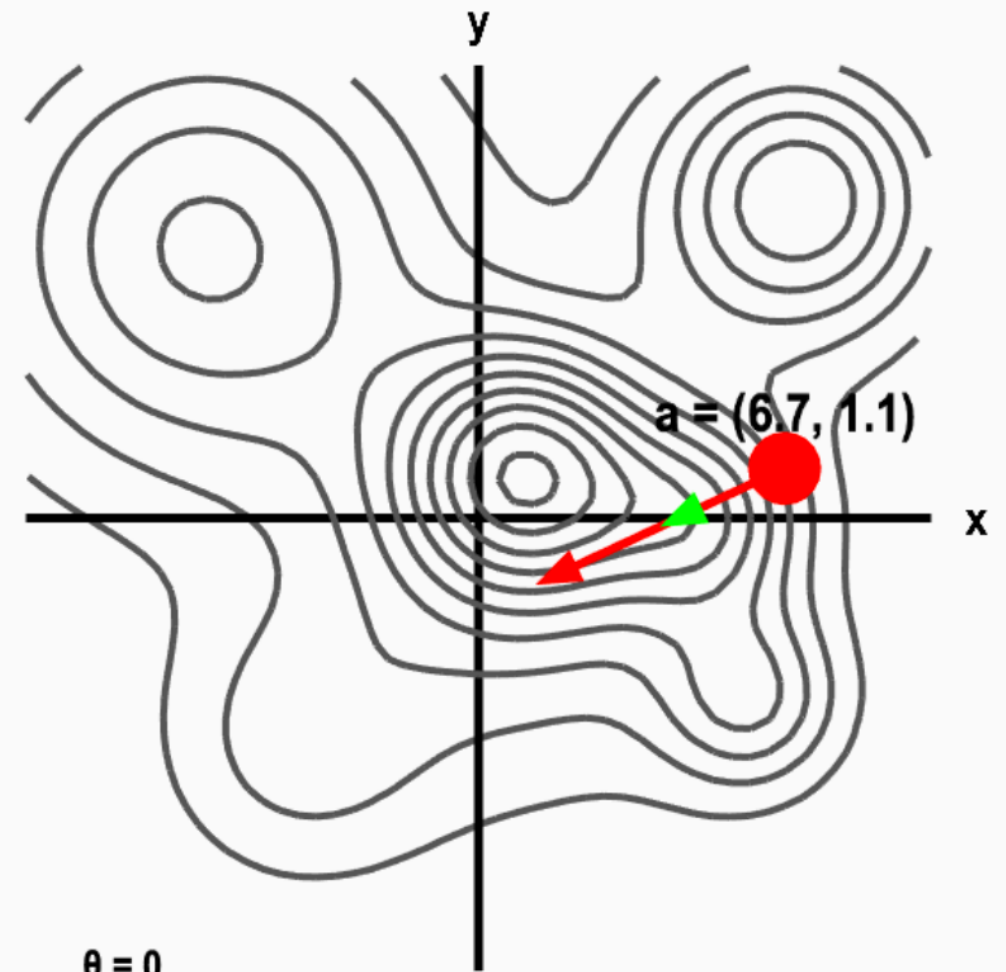
$u = (-0.91, -0.42)$

$a = (6.7, 1.1)$

$\nabla f(a) = (-1.81, -0.85)$

$D_u f(a) = 2.00$

$|\nabla f(a)| = 2.00$



$\theta = 0$



$u = (-0.91, -0.42)$

$\nabla f(a) = (-1.81, -0.85)$

$D_u f(a) = 2.00$

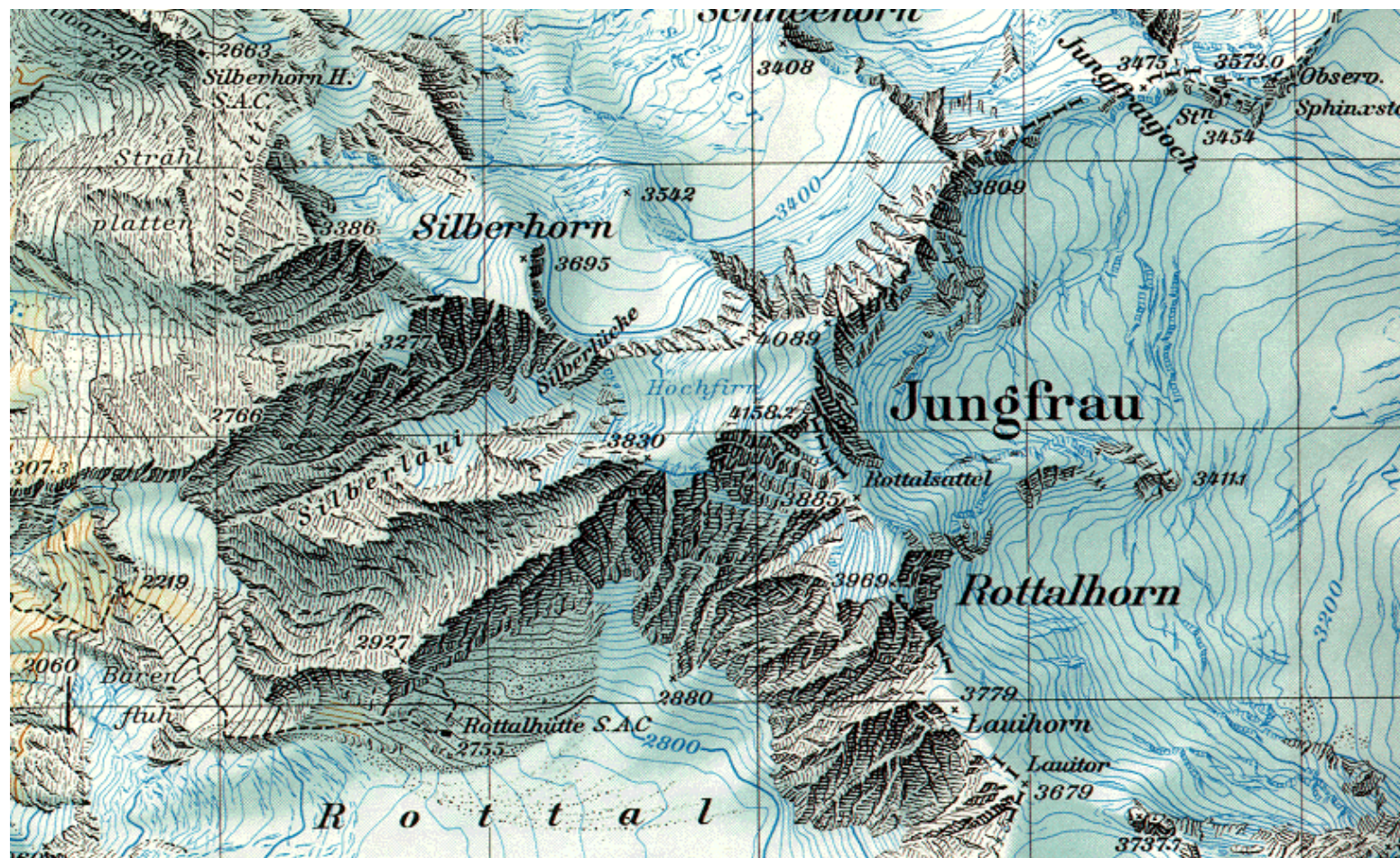
$|\nabla f(a)| = 2.00$

$f(a) = 4.87$

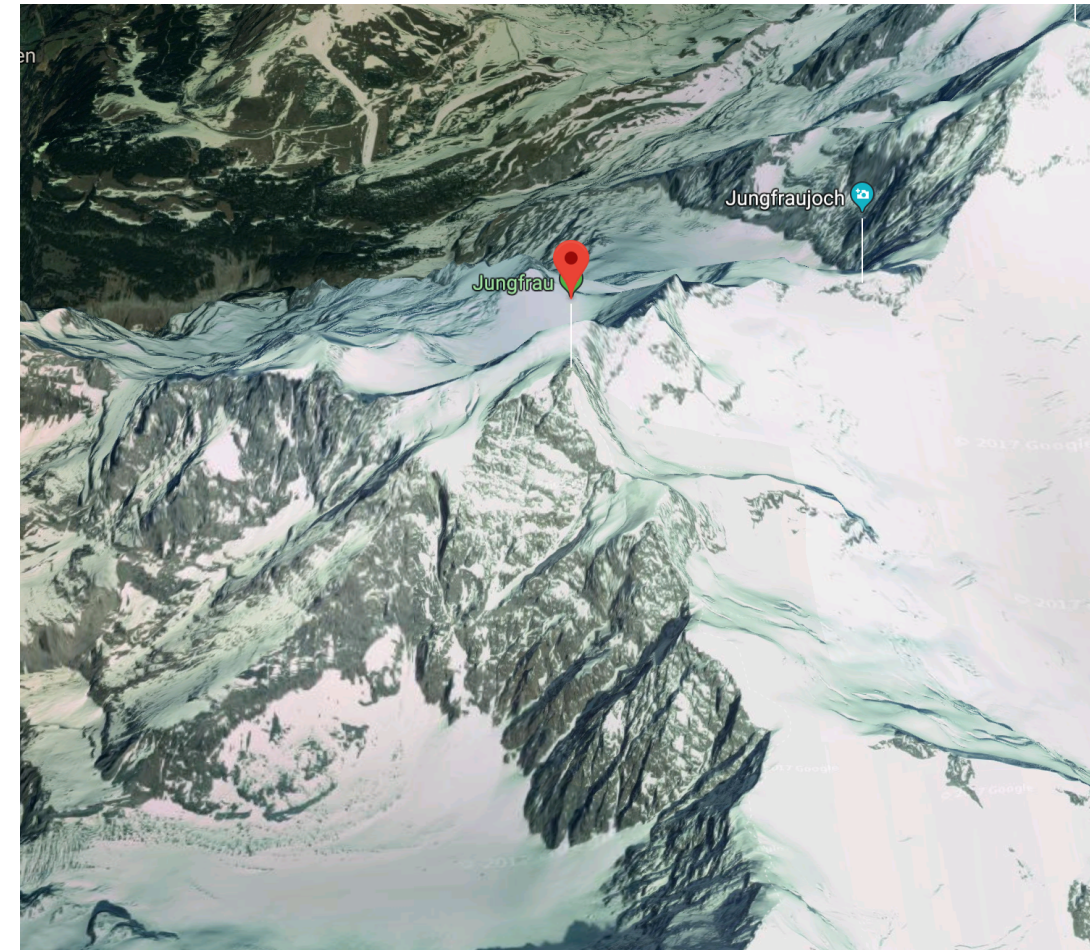


Level Sets: Topographic Map

The function is the altitude



Topographic map



3-D picture

Level Set: Exercise

Consider a strictly convex-quadratic function

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

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

$$f(x) = \frac{1}{2} (x - x^*)^T H (x - x^*) \quad H > 0$$

strictly convex.

$$f(x) \geq 0 \quad \text{because } H > 0$$

$$f(x) = 0 \Leftrightarrow x - x^* = 0 \Leftrightarrow x = x^*$$

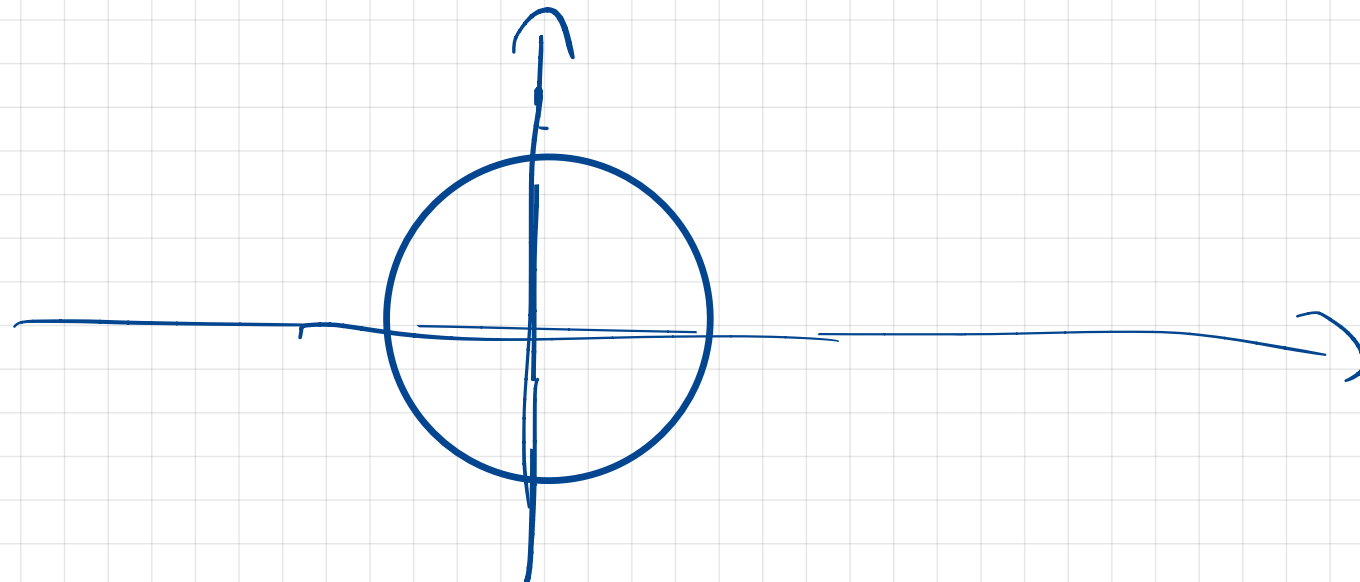
$$H = \nabla^2 f(x) \quad \text{Hessian matrix}$$

$$2/ \quad H = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\text{WLG } x^* = 0 \quad f(x) = \frac{1}{2} (x_1^2 + x_2^2)$$

$$\mathcal{L}_c = \left\{ x \mid \frac{1}{2} (x_1^2 + x_2^2) = c \right\}$$

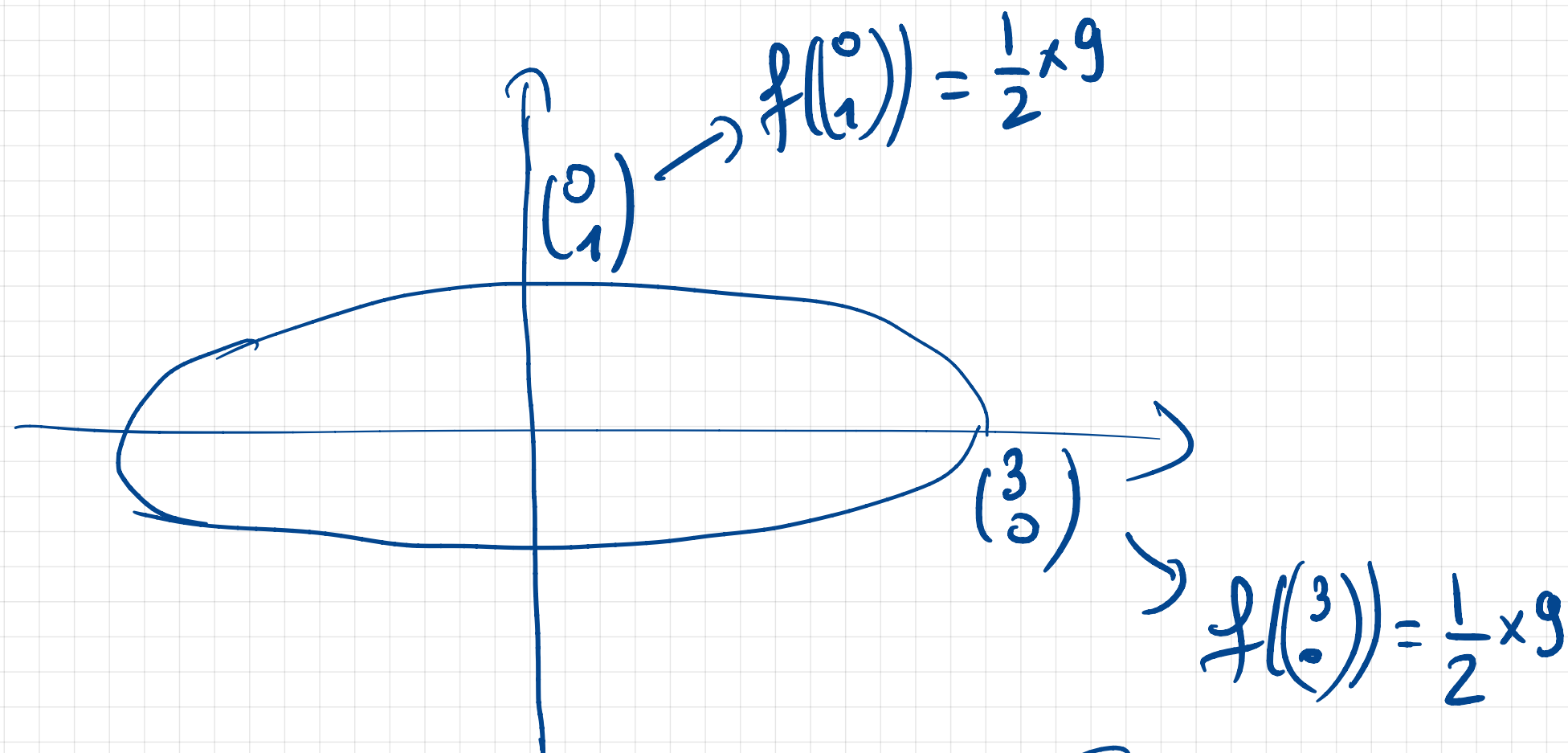
$$c > 0$$



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

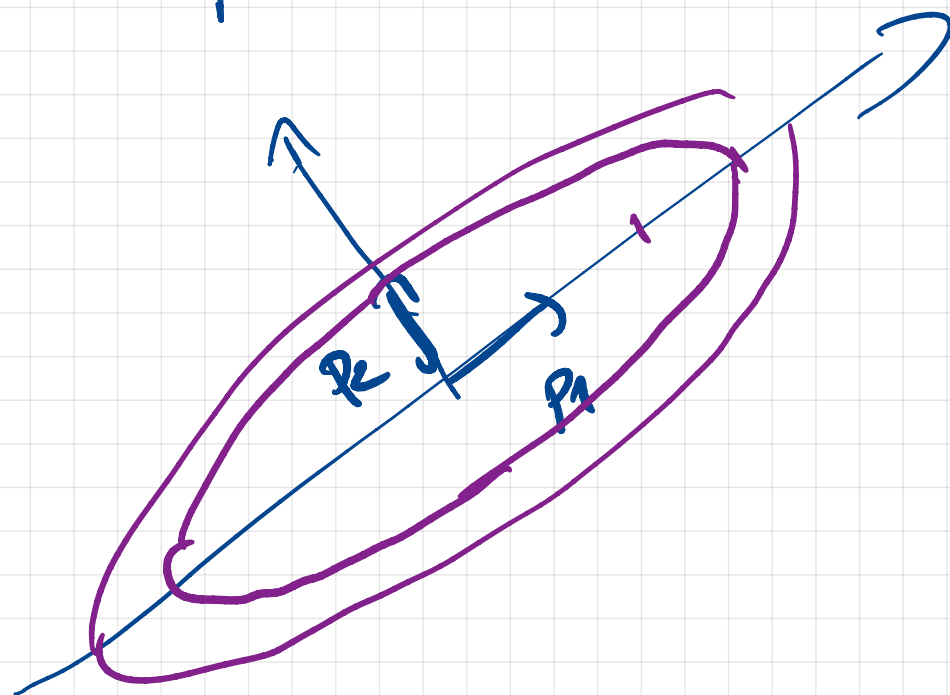
$$\mathcal{L}_c = \left\{ x = (x_1, x_2) \mid \frac{1}{2} (x_1^2 + 9x_2^2) = c \right\}$$

$c > 0$, ellipsoid.



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

$P = (p_1, p_2)$ eigenvectors of H

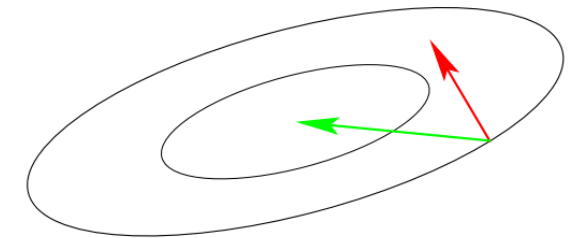
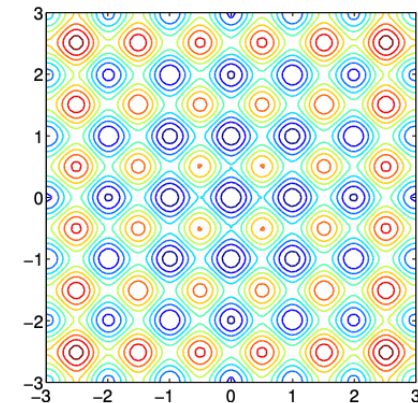
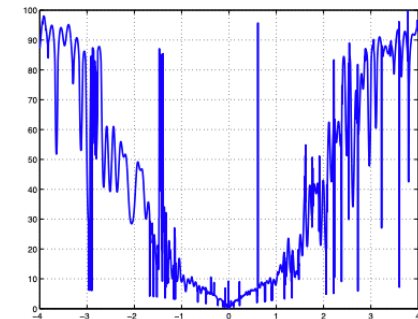
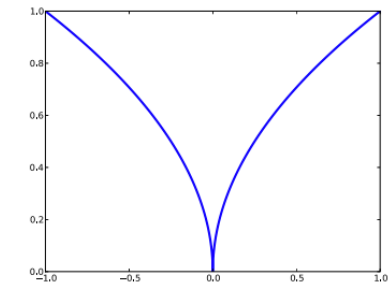


What Makes an Optimization Problem Difficult?

What Makes a Function Difficult to Solve?

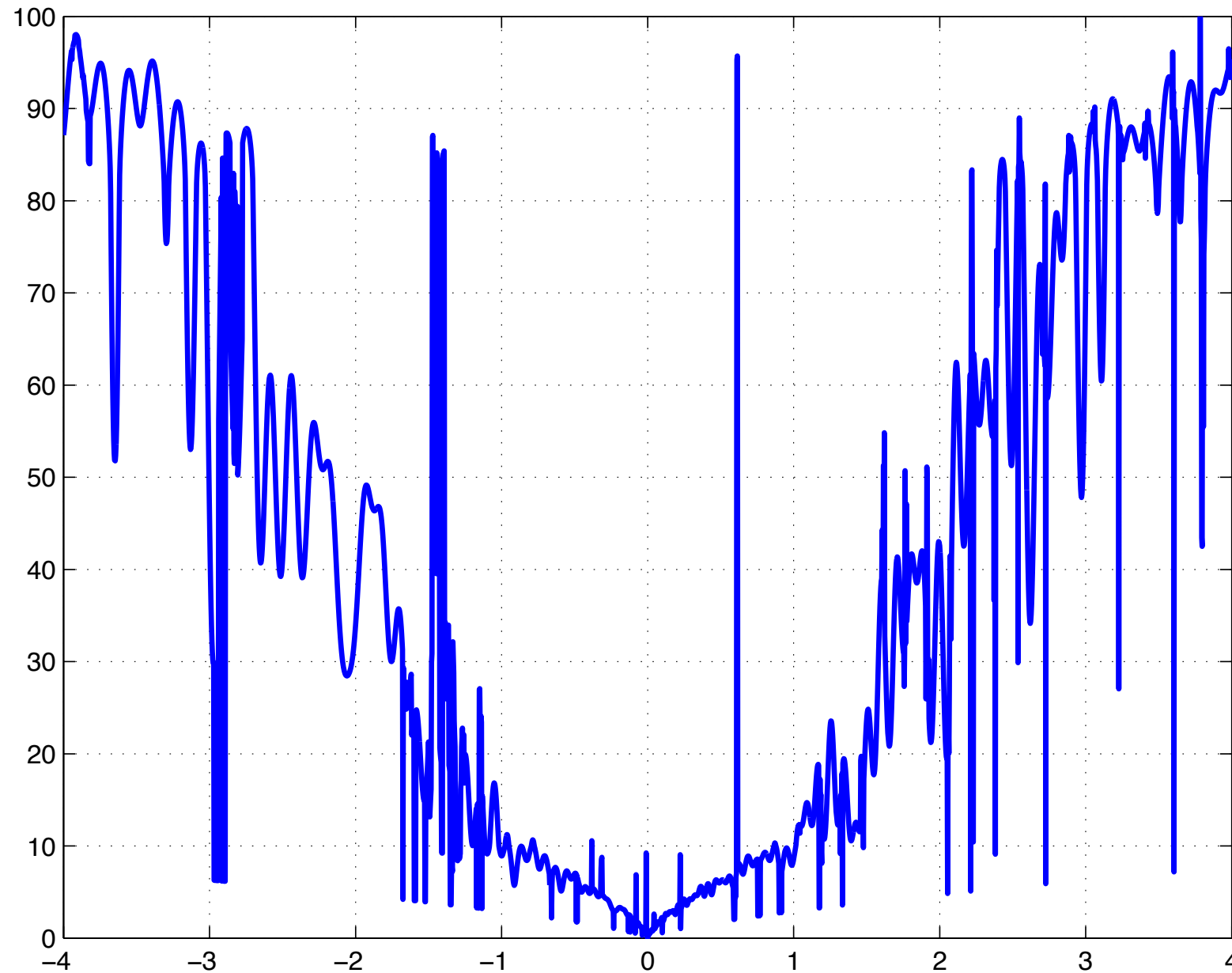
Why stochastic search?

- ▶ 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
- ▶ dimensionality (size of search space)
(considerably) larger than three
- ▶ non-separability
dependencies between the
objective variables
- ▶ ill-conditioning



gradient direction Newton direction

Ruggedness



$$f: x \in \mathbb{R}^4 \rightarrow \mathbb{R}$$

$$d \in \mathbb{R}^4$$

$$t \in \mathbb{R} \rightarrow f(td)$$

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

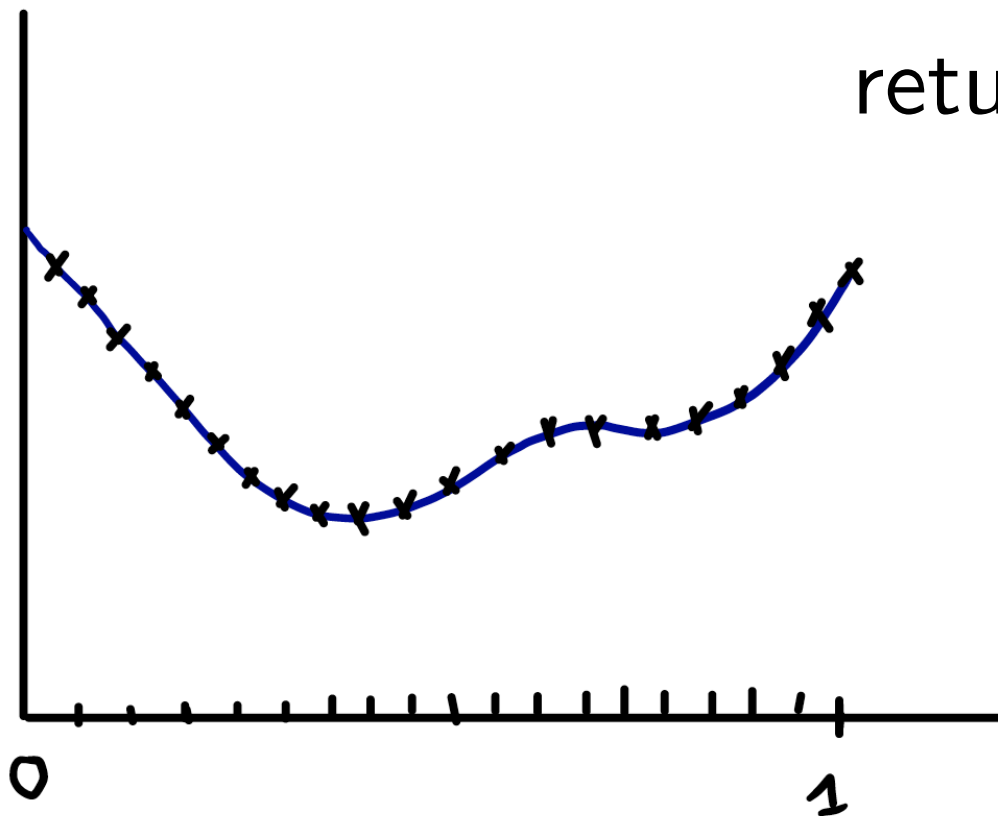
if $n=1$, which simple approach could you use to minimize:

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

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

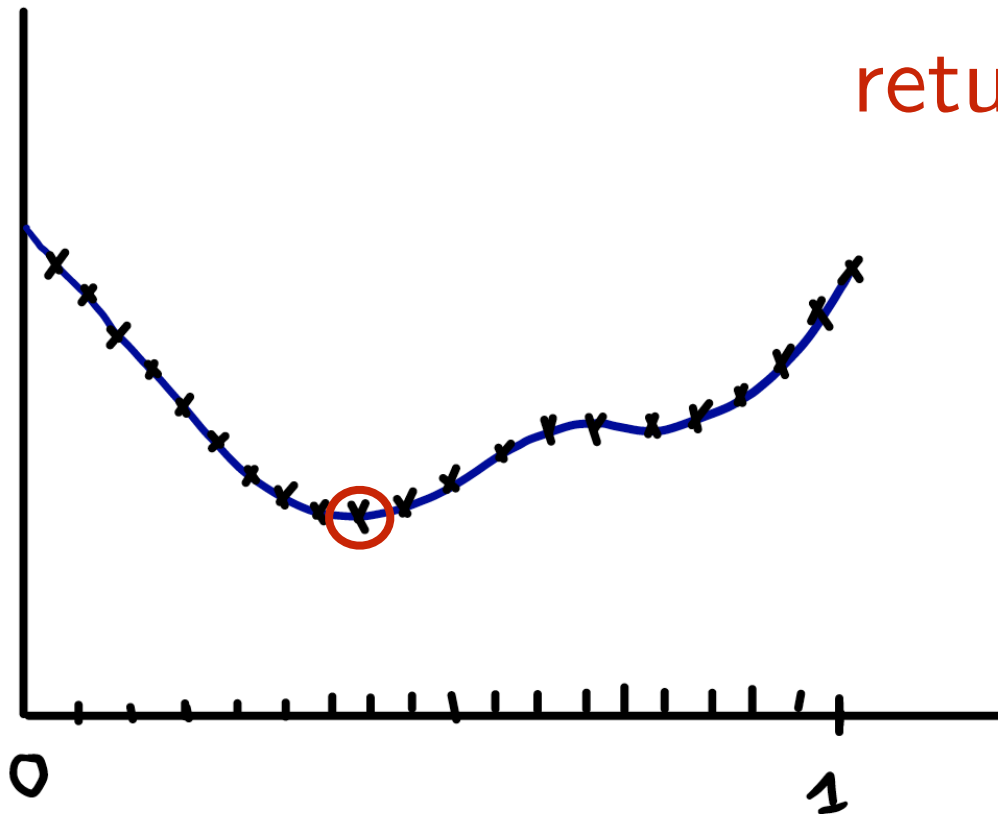
if $n=1$, which simple approach could you use to minimize:

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

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

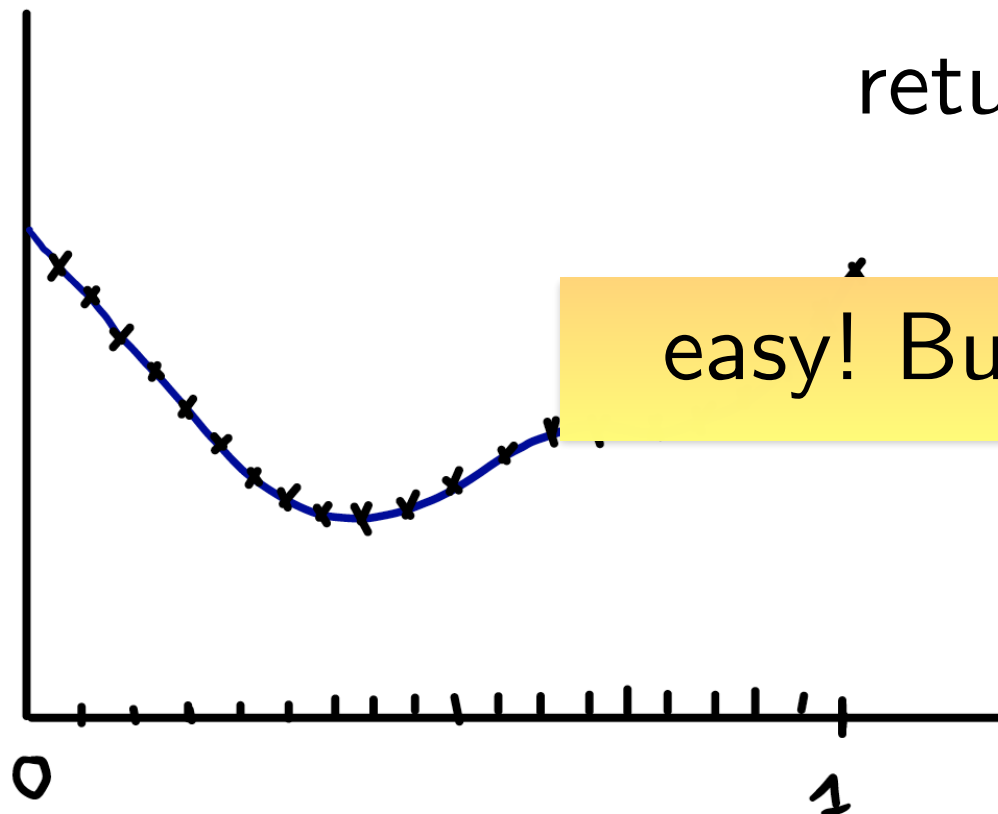
if $n=1$, which simple approach could you use to minimize:

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

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. ^{grid search} 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

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

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

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}

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

→ 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.

Proof: $\mathbb{R} \xrightarrow{h} \mathbb{R}$ $y \mapsto h(y)$ Let $g: \text{Im}(h) \rightarrow \mathbb{R}$ strict increasing

$$\underset{y}{\operatorname{argmin}} h(y) = \underset{y}{\operatorname{argmin}} g \circ h(y)$$

Let $\bar{x} \in \underset{y}{\operatorname{argmin}} h(y)$

$$h(\bar{x}) \leq h(y) \quad \forall y$$

Since g strict increasing

$$g \circ h(\bar{x}) \leq g \circ h(y) \quad \forall y$$

$$\Rightarrow \bar{x} \in \underset{y}{\operatorname{argmin}} g \circ h$$

$$\text{Let } \bar{x} \in \operatorname{argmin}_y g \circ h(y) \quad g \circ h(\bar{x}) \leq g \circ h(y) \quad \forall y$$

$$\xrightarrow[\text{generalized inverse}]{g^{-1}} g^{-1}(g \circ h(\bar{x})) \leq g^{-1}(g \circ h(y)) \quad \forall y$$

$$\Rightarrow h(\bar{x}) \leq h(y) \quad \forall y$$

$$\Rightarrow \bar{x} \in \operatorname{argmin} h$$

Since the argmin is preserved when composing with g strictly increasing to the left, then if f is separable, $g \circ f$ is separable.

Example

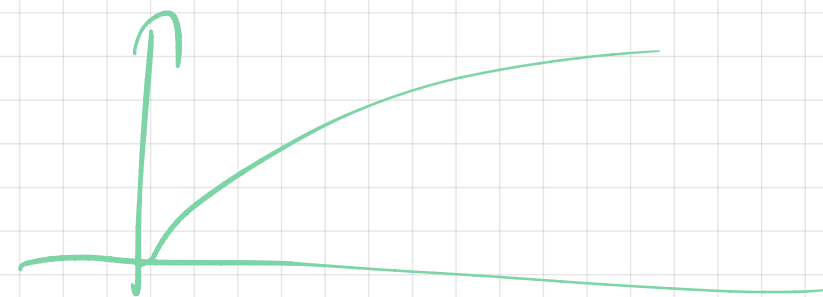
$$f(x) = \sum_{i=1}^n x_i^2$$

$$\tilde{f}(x) = \left(\sum_{i=1}^n x_i^2 \right)^{1/4}$$

"
 $g \circ f(x)$

$$g(x) = \begin{matrix} \mathbb{R}_{\geq 0} \\ \text{"} \\ \inf \\ \uparrow \\ x \end{matrix} \rightarrow \mathbb{R}$$

$$x \rightarrow x^{1/4}$$



$g \rightarrow$

$$\operatorname{argmin} h = \operatorname{argmax} g \circ f$$

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

Let us prove that $(\operatorname{argmin}_y f_{x_{\neg 1}}(y), \dots, \operatorname{argmin}_y f_{x_{\neg n}}^n(y)) \subset \operatorname{argmin} f$

$$x_i \in \operatorname{argmin}_{\alpha} f(x_1, \dots, x_{i-1}, \alpha, x_{i+1}, \dots, x_n) \quad i = 1, \dots, n$$

$$f(x) = f(x_1, \dots, x_n) \underset{\substack{\uparrow \\ \text{by def of } \alpha_1}}{\geq} f(\alpha_1, x_2, \dots, x_n) \underset{\substack{\uparrow \\ \text{by def of } \alpha_2}}{\geq} f(\alpha_1, \alpha_2, x_3, \dots, x_n)$$

$$f(x) \geq \dots \geq f(\alpha_1, \dots, \alpha_n) \quad \forall x$$

$$(\alpha_1, \dots, \alpha_n) \in \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} f.$$

The other inclusion is immediate:

$$\underset{x}{\operatorname{argmin}} f \subset (\underset{x}{\operatorname{argmin}} f(y), \dots,)$$

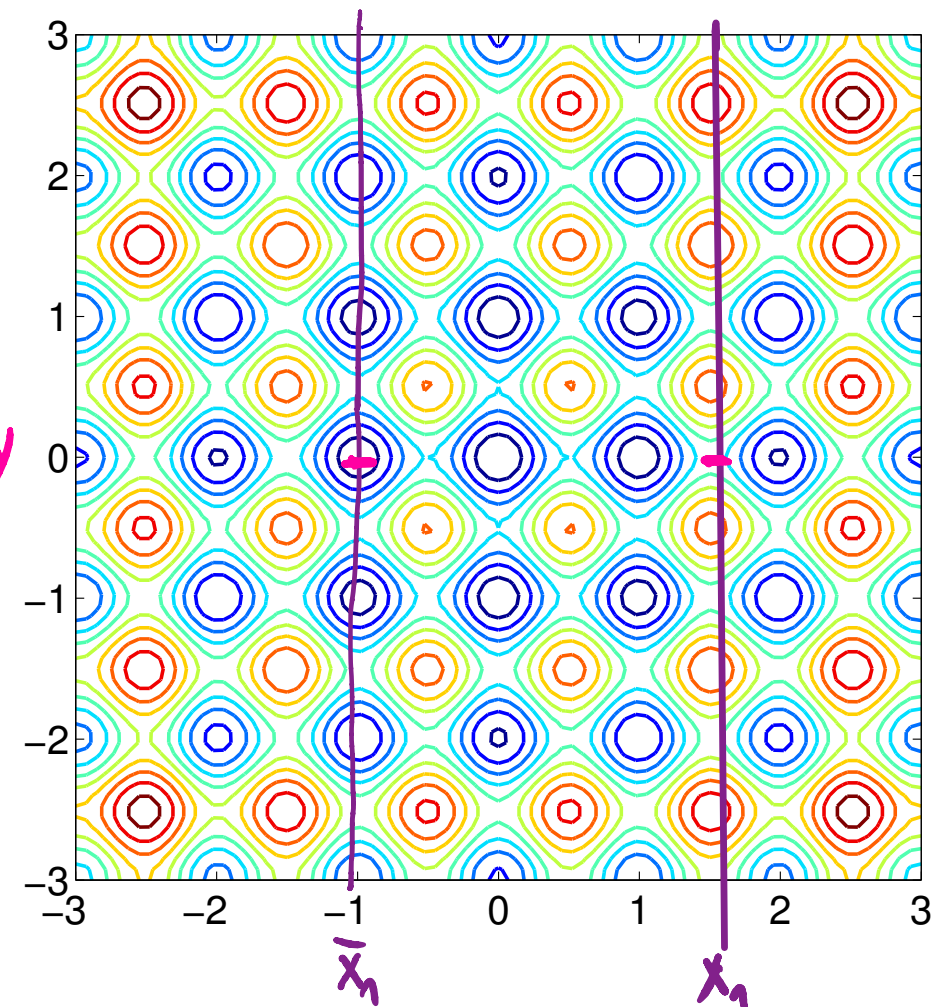
Example: Additively Decomposable Functions

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.

Proof:

$$f(x) = \exp \left(\ln \prod_{i=1}^n h_i(x_i) \right)$$

$$= \exp \left(\sum_{i=1}^n \ln h_i(x_i) \right)$$

$$= g \circ \text{additively decomposable}$$

$$g(x) = \exp(x) \text{ strict inc}$$

$$\hat{f}(x) = \sum_{i=1}^n \ln h_i(x_i) : \text{additively decomposable} \\ \hookrightarrow \text{separable.}$$

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)

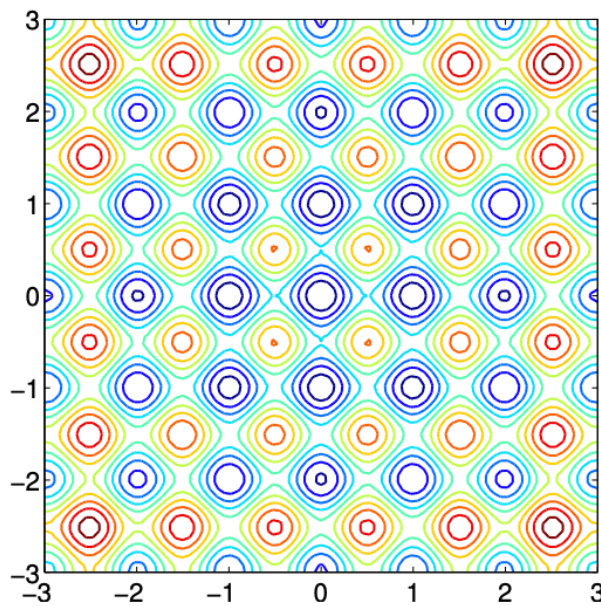
Non-separable Problems

Building a non-separable problem from a separable one

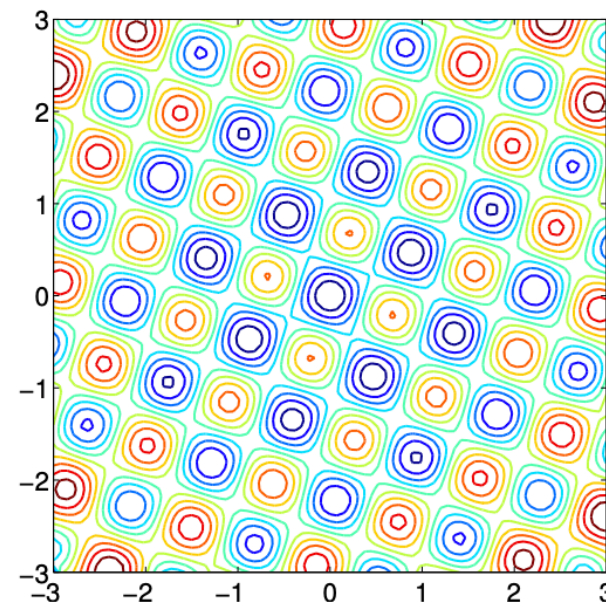
Rotating the coordinate system

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

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

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

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

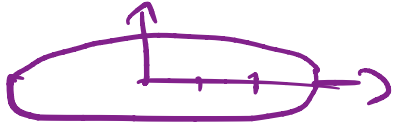
The condition number of the matrix H (with respect to the Euclidean norm) is defined as

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

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 $f(x) = \frac{1}{2}(x_1^2 + 9x_2^2)$

1. Compute its Hessian matrix, its condition number $H = \begin{pmatrix} 1 & 0 \\ 0 & 9 \end{pmatrix}$
 $\text{cond}(H) = 9$
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

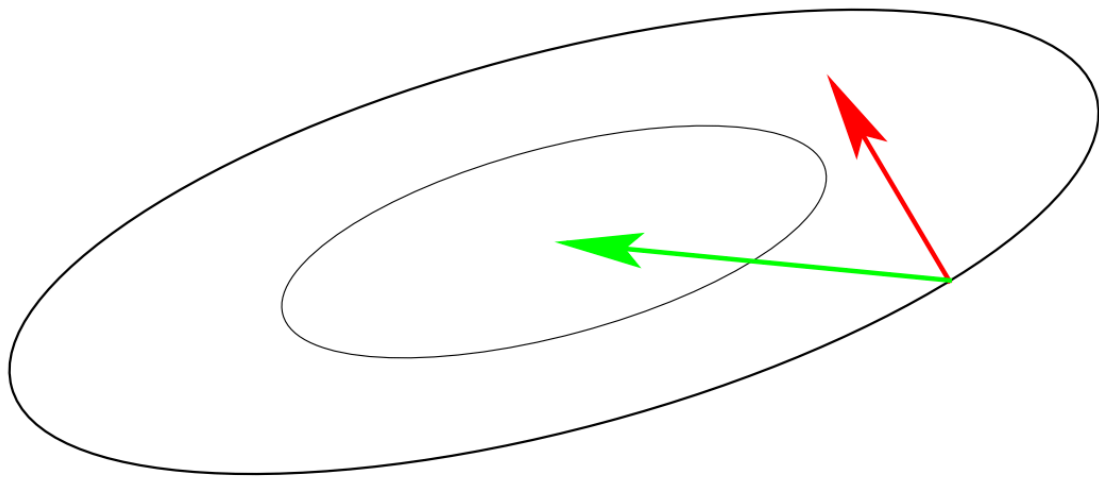
Real-world problems are often ill-conditioned.

4. Why do you think it is the case? \rightarrow physical variables optimized 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 $-f'(\mathbf{x})^T$

Newton direction

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

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

DERIVATIVE FREE OPTIMIZATION 2024/2025.

CLASS 2.

Pure Random Search (PRS)

Assume $f: [-1, 1]^n \rightarrow \mathbb{R}$
 $x = (x_1, \dots, x_n) \rightarrow f(x)$

PRS:

Initialize $x_{\text{best}} = \text{Unif}([-1, 1]^n)$ (uniform)

WHILE NOT HAPPY (while stop criterion not met)

Sample $x \sim \text{Unif}([-1, 1]^n)$

If $f(x) \leq f(x_{\text{best}})$

$x_{\text{best}} \leftarrow x$

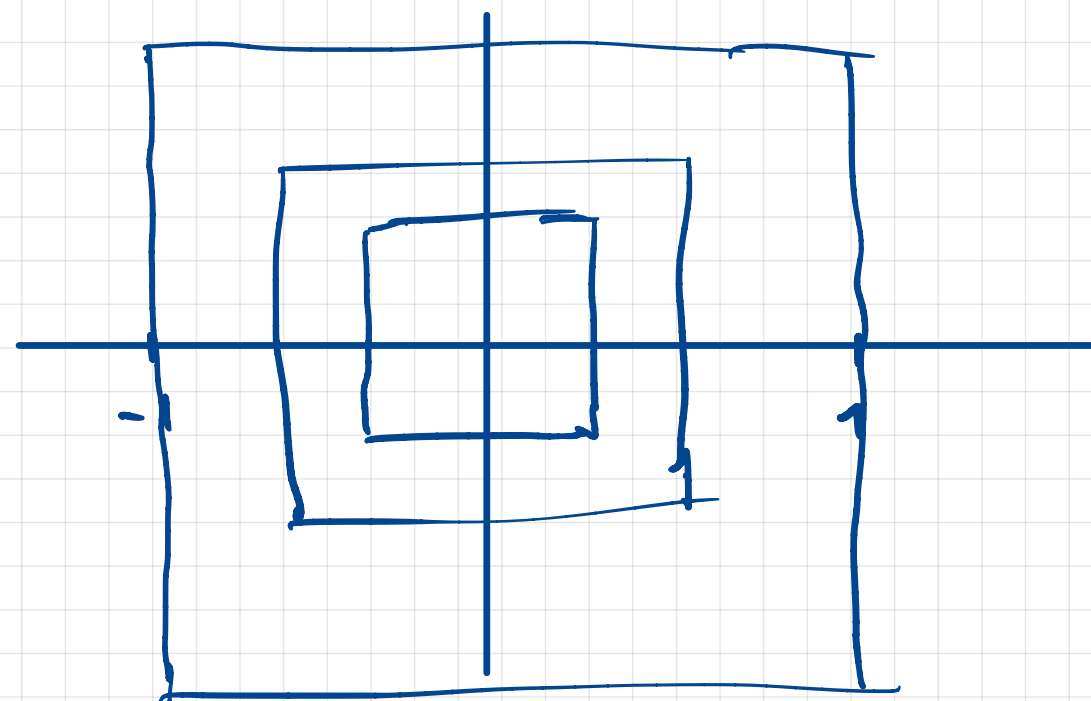
(see also Exercise 1)

Does this algorithm converge: Yes under mild assumptions on f
(need to have "volume" in a neighborhood of global optimum)

Simplified proof setting

$$f(x) = \|x\|_\infty = \max(|x_1|, \dots, |x_n|)$$

level sets : $n=2$



In exercise 1: uniform samples $\{u_0, u_1, \dots, u_t, \dots\}$
 $u_i \sim \text{Unif}([-1, 1]^n)$

x_t : best solution at iteration t

$$f(x_t) = \min \{f(u_1), \dots, f(u_t)\}$$

By induction.

Prove that $\forall \varepsilon > 0 \quad \lim_{t \rightarrow \infty} \mathbb{P}_t(\|X_t\|_\infty \geq \varepsilon) = 0$

↳ give CV in probability.

$$\|X_t\|_\infty = \min \{ \|V_1\|_\infty, \dots, \|V_t\|_\infty \}$$

$$\{ \|X_t\|_\infty \geq \varepsilon \} = \{ \forall_{k=1, \dots, t} \|V_k\|_\infty \geq \varepsilon \} = \bigcap_{k=1}^t \{ \|V_k\|_\infty \geq \varepsilon \}$$

$$\mathbb{P}_r(\{ \|X_t\|_\infty \geq \varepsilon \}) = \mathbb{P}_r\left(\bigcap_{k=1}^t \{ \|V_k\|_\infty \geq \varepsilon \}\right)$$

$$= \prod_{k=1}^t \mathbb{P}_r(\|V_k\|_\infty \geq \varepsilon) = \mathbb{P}_r(\|V_1\|_\infty \geq \varepsilon)^t$$

by ind of $\{V_k, k=1, \dots, t\}$

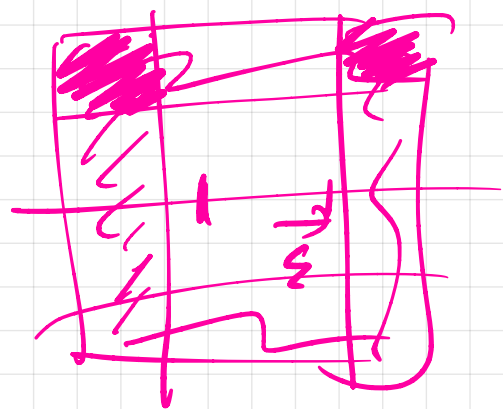
↑
because V_k are identically distributed.
 $= (1 - \varepsilon^n)^t$

$$\Pr(\|U_1\|_\infty > \epsilon) \stackrel{?}{=} 1 - \Pr(\|U_1\|_\infty \leq \epsilon)$$

X

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

$$\prod_{\text{Coord}} (\text{each coordinate } > \epsilon) = 1 - \frac{(2\epsilon)^n}{2^n} = 1 - \epsilon^n$$

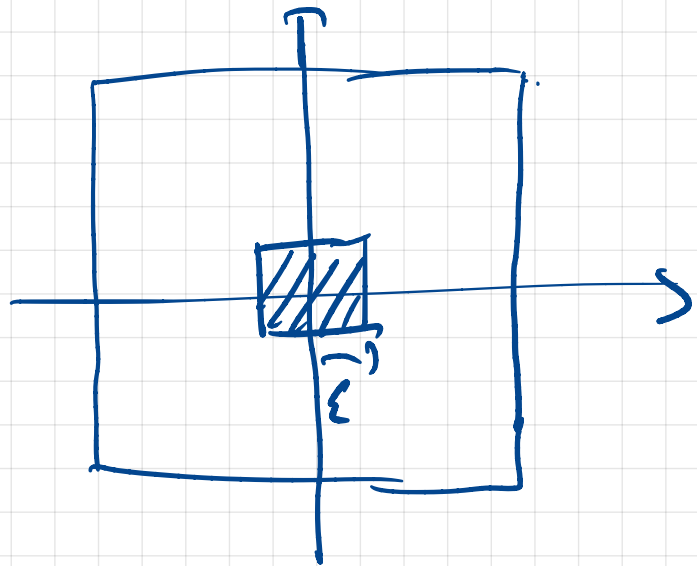


$$\Pr(\|X_t\|_\infty > \epsilon) = (1 - \epsilon^n)^t \xrightarrow{t \rightarrow \infty} 0$$

Hence IRS converges in probability to the optimum of f .

Let's look at how fast it converges.

$$T_\epsilon = \inf \{ t \mid X_t \in \underbrace{B(x^*, \epsilon)}_{\{x \mid \|x\|_\infty \leq \epsilon\}} \} \quad (\text{Ball for infinity norm})$$



PRS is similar to game:

Sample U_t , win if $f(U_t) \leq \epsilon$
 $(\Leftrightarrow) U_t \in B(0, \epsilon)$

Loose otherwise.

T_ϵ : time it takes to win this game.

Given a game with 2 outcomes win with proba p and loose with proba $(1-p)$, where the outcome is sampled randomly and independently [example: flip a coin], the time it takes to win a game is distributed according to a geometric distribution.

$$\mathbb{E}[T_\epsilon] = \frac{1}{p}$$

Back to PRS. $p = \Pr(\text{"win"}) = \Pr(\|U_t\| \leq \epsilon) = \epsilon^n$

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

Is it fast? No

To compare: linear cv $T_\varepsilon \sim n \log\left(\frac{1}{\varepsilon}\right)$ [gradient des on L_2 norm strongly conv]

The algorithm is "blind", does not take into account the information gathered on f , through the sampling of points to sample "better" solutions.

Part II: Algorithms

Landscape of Derivative Free Optimization Algorithms

Deterministic Algorithms

Quasi-Newton with estimation of gradient (BFGS) [Broyden et al. 1970]

Simplex downhill [Nelder and Mead 1965]

Pattern search, Direct Search [Hooke and Jeeves 1961]

Trust-region/Model Based methods (NEWUOA, BOBYQA) [Powell, 06,09]

Stochastic (randomized) search methods

Evolutionary Algorithms (continuous domain)

Differential Evolution [Storn, Price 1997]

Particle Swarm Optimization [Kennedy and Eberhart 1995]

Evolution Strategies, CMA-ES [Rechenberg 1965, Hansen, Ostermeier 2001]

Estimation of Distribution Algorithms (EDAs) [Larrañaga, Lozano, 2002]

Cross Entropy Method (same as EDAs) [Rubinstein, Kroese, 2004]

Genetic Algorithms [Holland 1975, Goldberg 1989]

Simulated Annealing [Kirkpatrick et al. 1983]

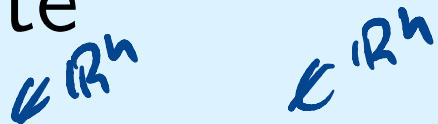
A Generic Template for Stochastic Search

Define $\{P_\theta : \theta \in \Theta\}$, a family of probability distributions on \mathbb{R}^n

Generic template to optimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

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

While not terminate

1. Sample x_1, \dots, x_λ according to P_θ

2. Evaluate x_1, \dots, x_λ on f
3. Update parameters $\theta \leftarrow F(\theta, x_1, \dots, x_\lambda, f(x_1), \dots, f(x_\lambda))$

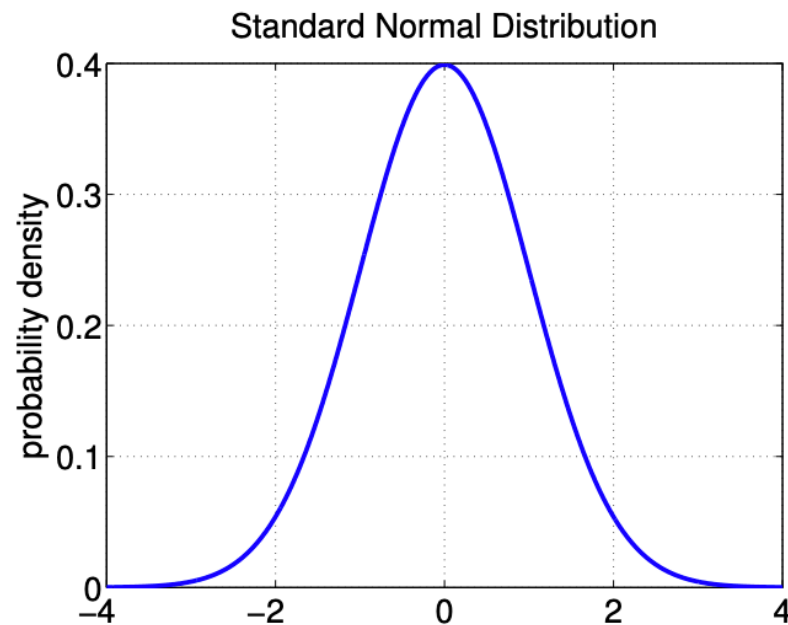
the update of θ should drive P_θ to concentrate on the optima of f

To obtain an optimization algorithm we need:

- ❶ to define $\{P_\theta, \theta \in \Theta\}$
- ❷ to define F the update function of θ

Which probability distribution to sample candidate solutions?

Normal distribution - 1D case



probability density of the 1-D standard normal distribution $\mathcal{N}(0, 1)$

(expected (mean) value, variance) = (0,1)

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

General case

► Normal distribution $\mathcal{N}(\mathbf{m}, \sigma^2)$

(expected value, variance) = (\mathbf{m}, σ^2)

density: $p_{\mathbf{m},\sigma}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mathbf{m})^2}{2\sigma^2}\right)$

- A normal distribution is entirely determined by its mean value and variance
- The family of normal distributions is closed under linear transformations: if X is normally distributed then a linear transformation $aX + b$ is also normally distributed
- **Exercise:** Show that $\mathbf{m} + \sigma\mathcal{N}(0, 1) = \mathcal{N}(\mathbf{m}, \sigma^2)$

$m + \sigma \mathcal{N}(0,1)$ is normally distributed

We only need to identify its mean and variance:

$$\mathbb{E}(m + \sigma \mathcal{N}(0,1)) \underset{\substack{\text{by linearity} \\ \text{on } \mathbb{E}}}{=} m + \sigma \underbrace{\mathbb{E}(\mathcal{N}(0,1))}_{=0} = m$$

$$\begin{aligned} & \text{Var}(m + \sigma \mathcal{N}(0,1)) \\ &= \mathbb{E}\left(\left(m + \sigma \mathcal{N}(0,1) - m\right)^2\right) - \\ &= \mathbb{E}\left(\sigma^2 \mathcal{N}(0,1)^2\right) \\ &= \sigma^2 \underbrace{\mathbb{E}(\mathcal{N}(0,1)^2)}_{=1} = \sigma^2 \end{aligned}$$

$$\Rightarrow m + \sigma \mathcal{N}(0,1) \approx \mathcal{N}(m, \sigma^2)$$

$$\mathbb{P}(m + \sigma \mathcal{N}(0,1) \leq t) = \mathbb{P}_r \left(\mathcal{N}(0,1) \leq \frac{t-m}{\sigma} \right)$$

$$= \int_{-\infty}^{\frac{t-m}{\sigma}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) dx$$

$$dy = \sigma dx$$

$$y = \sigma x + m$$

$$= \int_{-\infty}^t \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\left(\frac{y-m}{\sigma}\right)^2\right) dy$$

$$x = \frac{t-m}{\sigma} \quad y = t$$

Generalization to n Variables: Independent Case

Assume $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ denote its density $p(x_1) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2\right)$

Assume $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ denote its density $p(x_2) = \frac{1}{Z_2} \exp\left(-\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2\right)$

Assume X_1 and X_2 are **independent**, then (X_1, X_2) is a Gaussian vector with

$$p(x_1, x_2) =$$

Generalization to n Variables: Independent Case

Assume $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ denote its density $p(x_1) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2\right)$

Assume $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ denote its density $p(x_2) = \frac{1}{Z_2} \exp\left(-\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2\right)$

Assume X_1 and X_2 are **independent**, then (X_1, X_2) is a Gaussian vector with

$$p(x_1, x_2) = p(x_1)p(x_2) = \frac{1}{Z_1 Z_2} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$$\text{with } x = (x_1, x_2)^T \quad \mu = (\mu_1, \mu_2)^T \quad \Sigma = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$$

Generalization to n Variables: Independent Case

Assume $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ denote its density $p(x_1) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2\right)$

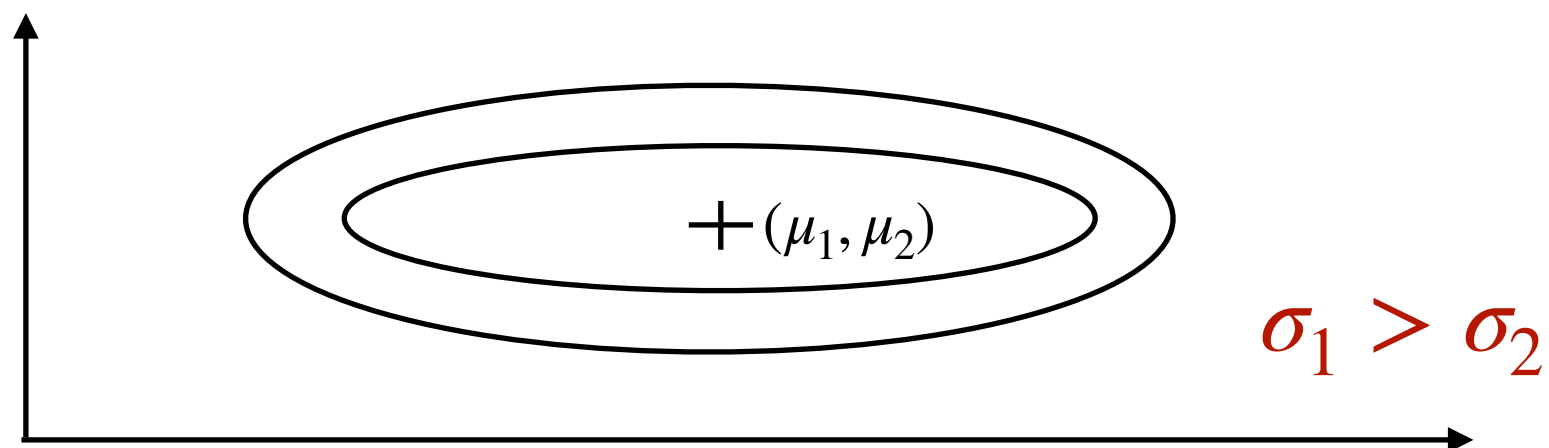
Assume $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ denote its density $p(x_2) = \frac{1}{Z_2} \exp\left(-\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2\right)$

Assume X_1 and X_2 are **independent**, then (X_1, X_2) is a Gaussian vector with

$$p(x_1, x_2) = p(x_1)p(x_2) = \frac{1}{Z_1 Z_2} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$\rightarrow -\frac{1}{2} \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right]$

with $x = (x_1, x_2)^T$ $\mu = (\mu_1, \mu_2)^T$ $\Sigma = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$



Generalization to n Variables: General Case

Gaussian Vector - Multivariate Normal Distribution

A random vector $X = (X_1, \dots, X_n) \in \mathbb{R}^n$ is a **Gaussian vector** (or multivariate normal) if and only if for all real numbers a_1, \dots, a_n , the random variable $a_1X_1 + \dots + a_nX_n$ has a **normal distribution**.

Gaussian Vector - Multivariate Normal Distribution

A random variable following a 1-D normal distribution is determined by its mean value m and variance σ^2 .

In the n -dimensional case it is determined by its mean vector and covariance matrix

Covariance Matrix

If the entries in a vector $\mathbf{X} = (X_1, \dots, X_n)^T$ are random variables, each with finite variance, then the covariance matrix Σ is the matrix whose (i, j) entries are the covariance of (X_i, X_j)

$$\Sigma_{ij} = \text{cov}(X_i, X_j) = \mathbb{E} [(X_i - \mu_i)(X_j - \mu_j)]$$

where $\mu_i = \mathbb{E}(X_i)$. Considering the expectation of a matrix as the expectation of each entry, we have

$$\Sigma_{ii} = \text{Var}(X_i)$$

$$\Sigma = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T]$$

Σ is symmetric, positive definite

Density of a n-dimensional Gaussian vector $\mathcal{N}(m, C)$:

$$p_{\mathcal{N}(m, C)}(x) = \frac{1}{(2\pi)^{n/2} |C|^{1/2}} \exp \left(-\frac{1}{2} (x - m)^\top C^{-1} (x - m) \right)$$

$$|C| = \det(C)$$

The **mean vector** m :

determines the displacement

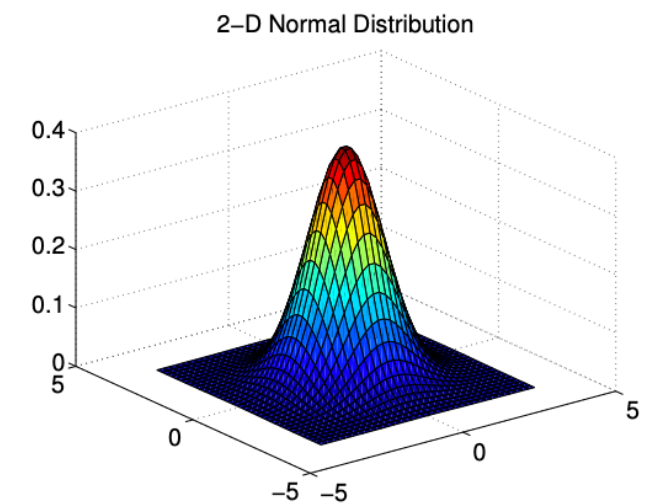
is the value with the largest density

the distribution is symmetric around the mean

$$\mathcal{N}(m, C) = m + \mathcal{N}(0, C)$$

The **covariance matrix**:

determines the geometrical shape (see next slides)



Geometry of a Gaussian Vector

Consider a Gaussian vector $\mathcal{N}(m, C)$, remind that lines of equal densities are given by:

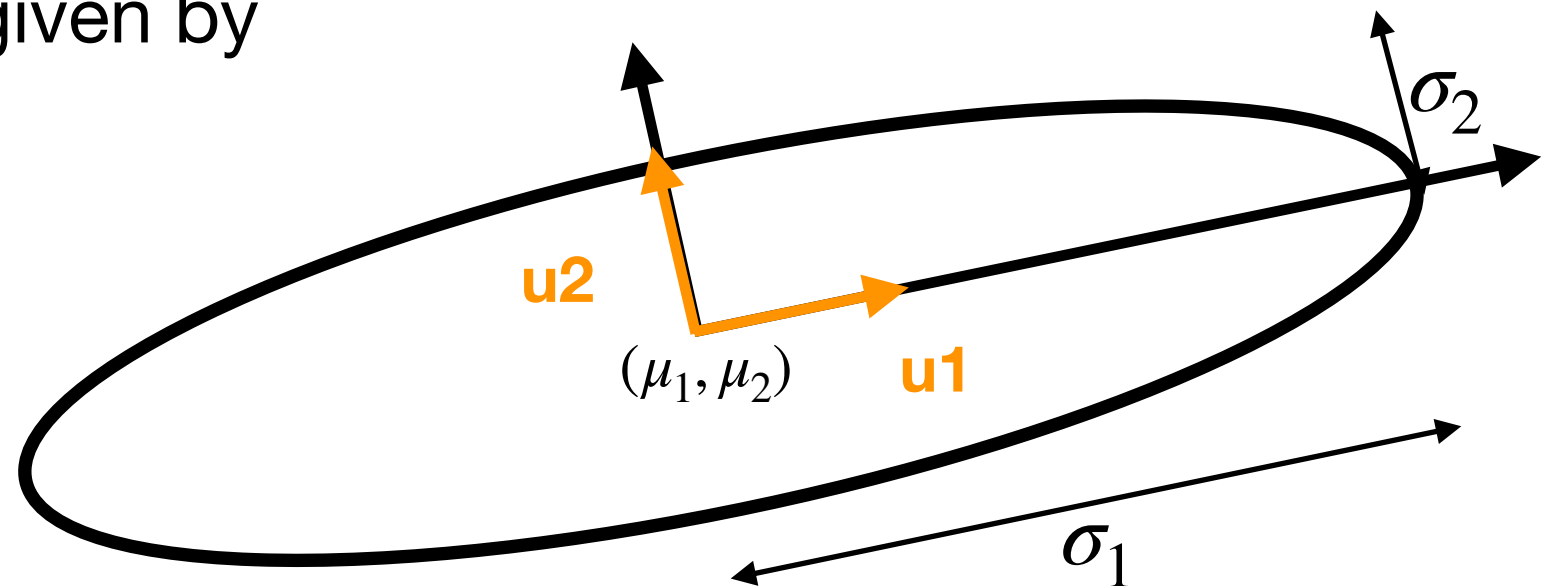
$$\{x \mid \Delta^2 = (x - m)^T C^{-1} (x - m) = \text{cst}\}$$

Decompose $C = U\Lambda U^T$ with U orthogonal, i.e.

$$C = \begin{pmatrix} u_1 & u_2 \\ | & | \end{pmatrix} \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix} \begin{pmatrix} u_1 & - \\ u_2 & - \end{pmatrix}$$

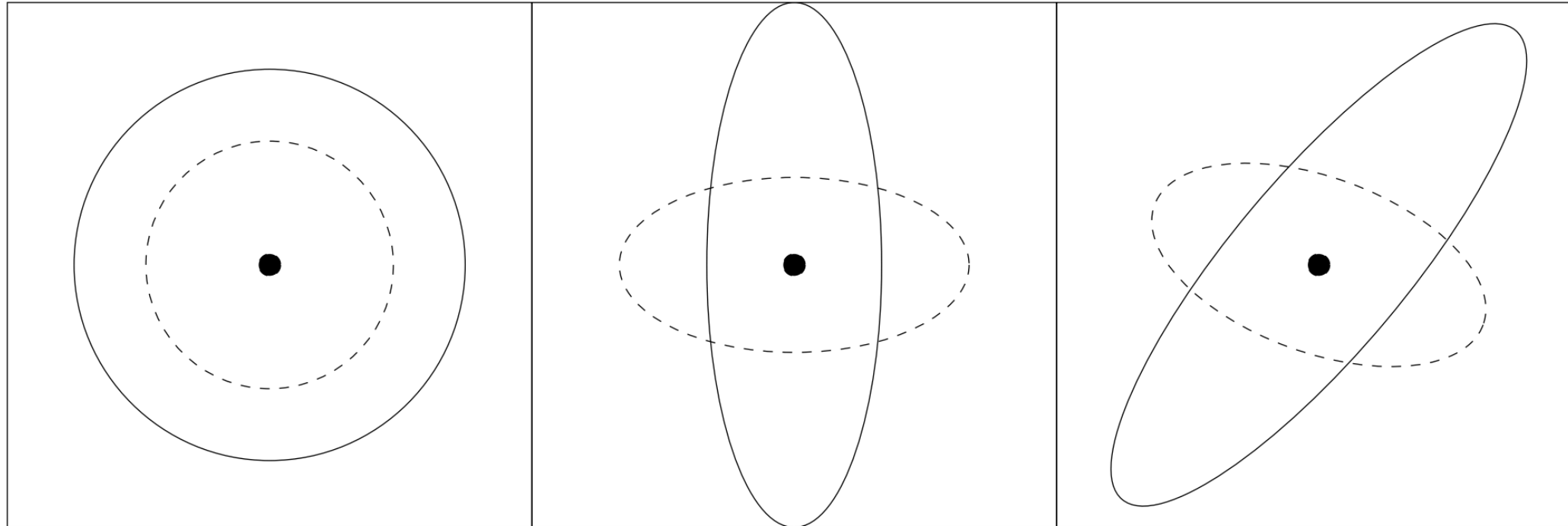
Let $Y = U^T(x - m)$, then in the coordinate system, (u_1, u_2) , the lines of equal densities are given by

$$\{x \mid \Delta^2 = \frac{Y_1^2}{\sigma_1^2} + \frac{Y_2^2}{\sigma_2^2} = \text{cst}\}$$



... any **covariance matrix** can be uniquely identified with the iso-density ellipsoid $\{\mathbf{x} \in \mathbb{R}^n \mid (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}) = 1\}$

Lines of Equal Density



$\mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{I}) \sim \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{I})$
 one degree of freedom σ
 components are
 independent standard
 normally distributed

$\mathcal{N}(\mathbf{m}, \mathbf{D}^2) \sim \mathbf{m} + \mathbf{D} \mathcal{N}(\mathbf{0}, \mathbf{I})$
 n degrees of freedom
 components are
 independent, scaled

$\mathcal{N}(\mathbf{m}, \mathbf{C}) \sim \mathbf{m} + \mathbf{C}^{\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I})$
 $(n^2 + n)/2$ degrees of freedom
 components are
 correlated

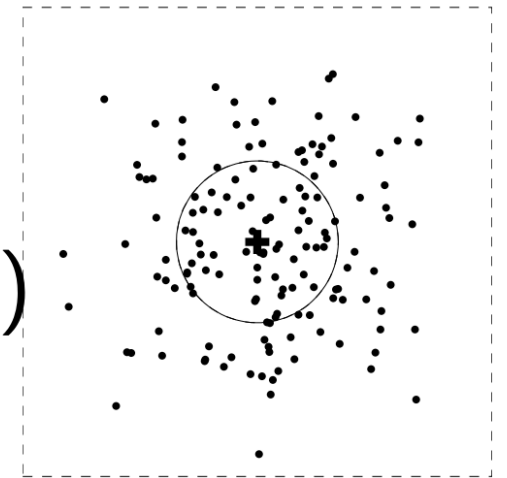
where \mathbf{I} is the identity matrix (isotropic case) and \mathbf{D} is a diagonal matrix (reasonable for separable problems) and $\mathbf{A} \times \mathcal{N}(\mathbf{0}, \mathbf{I}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^T)$ holds for all \mathbf{A} .

Evolution Strategies

New search points are sampled normally distributed

$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i$ for $i = 1, \dots, \lambda$ with \mathbf{y}_i i.i.d. $\sim \mathcal{N}(\mathbf{0}, \mathbf{C})$:

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

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{C})$$

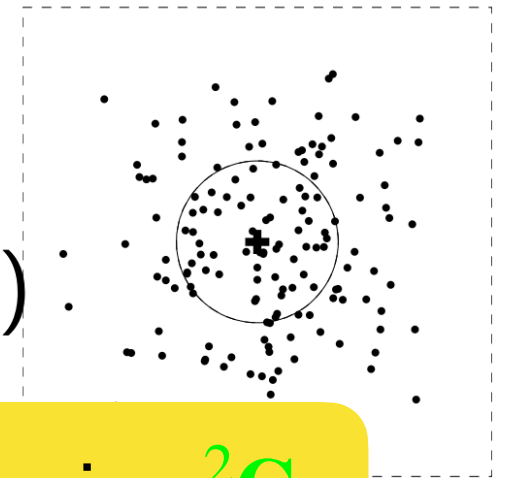
- ▶ the **mean** vector $\mathbf{m} \in \mathbb{R}^n$ represents the favorite solution
- ▶ the so-called **step-size** $\sigma \in \mathbb{R}_+$ controls the *step length*
- ▶ the **covariance matrix** $\mathbf{C} \in \mathbb{R}^{n \times n}$ determines the **shape** of the distribution ellipsoid

here, all new points are sampled with the same parameters

Evolution Strategies

New search points are sampled normally distributed

$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i$ for $i = 1, \dots, \lambda$ with \mathbf{y}_i i.i.d. $\sim \mathcal{N}(\mathbf{0}, \mathbf{C})$:



In fact, the covariance matrix of the sampling distribution is $\sigma^2 \mathbf{C}$ but it is convenient to refer to \mathbf{C} as the covariance matrix (it is a covariance matrix but not of the sampling distribution)

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

How to update the different parameters m, σ, \mathbf{C} ?

- 1. Adapting the mean m**
2. Adapting the step-size σ
3. Adapting the covariance matrix \mathbf{C}

Update the Mean: a Simple Algorithm the (1+1)-ES

Notation and Terminology:

one solution kept
from one iteration
to the next

(**1**+**1**)-ES

one new solution
(offspring) sampled at
each iteration

The $+$ means that we keep the best between current solution and new solution, we talk about *elitist selection*

(1+1)-ES algorithm (update of the mean)

sample one candidate solution from the mean \mathbf{m}

$$\mathbf{x} = \mathbf{m} + \sigma \mathcal{N}(0, \mathbf{C})$$

if \mathbf{x} is better than \mathbf{m} (i.e. if $f(\mathbf{x}) \leq f(\mathbf{m})$), select \mathbf{m}

$$\mathbf{m} \leftarrow \mathbf{x}$$

The (1+1)-ES algorithm is a simple algorithm, yet:

- the elitist selection is not robust to outliers

we cannot lose solutions accepted by “chance”, for instance that look good because the noise gave it a low function value

- there is no population (just a single solution is sampled) which makes it less robust

In practice, one should rather use a:

$(\mu/\mu, \lambda)$ -ES

The μ best solutions are
selected and recombined
(to form the new mean)

λ solutions are
sampled
at each iteration

The $(\mu/\mu, \lambda)$ -ES - Update of the Mean Vector

Given the i -th solution point $\mathbf{x}_i = \mathbf{m} + \sigma \underbrace{\mathbf{y}_i}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{C})} \quad i=1, \dots, \lambda$

Let $\mathbf{x}_{i:\lambda}$ the i -th ranked solution point, such that $f(\mathbf{x}_{1:\lambda}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda})$.

Notation: we denote $\mathbf{y}_{i:\lambda}$ the vector such that $\mathbf{x}_{i:\lambda} = \mathbf{m} + \sigma \mathbf{y}_{i:\lambda}$

Exercise: realize that $\mathbf{y}_{i:\lambda}$ is generally not distributed as $\mathcal{N}(\mathbf{0}, \mathbf{C})$

The new mean reads

$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}$$

where

$$w_1 \geq \dots \geq w_{\mu} > 0, \quad \sum_{i=1}^{\mu} w_i = 1, \quad \frac{1}{\sum_{i=1}^{\mu} w_i^2} =: \mu_w \approx \frac{\lambda}{4} \quad (\text{typically } \mu = \frac{\lambda}{2})$$

The best μ points are selected from the new solutions (non-elitistic) and weighted intermediate recombination is applied.

What changes in the previous slide if instead of optimizing f , we optimize $g \circ f$ where $g : \text{Im}(f) \rightarrow \mathbb{R}$ is strictly increasing?

$$f(x) = x^2$$

$$g \circ f(x) = (x^2)^{1/4} \\ = \sqrt{|x|}$$

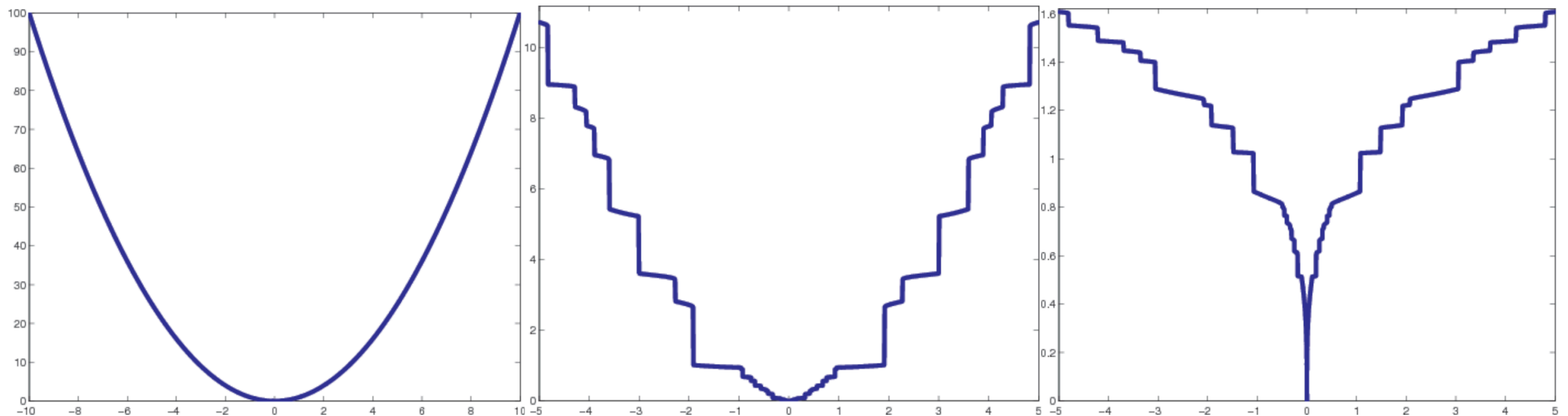
$$g(x) = x \in \mathbb{R}^+ \rightarrow x^{1/4}$$

Invariance Under Monotonically Increasing Functions

Comparison-based/ranking-based algorithms:

Update of all parameters uses only the ranking:

$$f(x_{1:\lambda}) \leq f(x_{2:\lambda}) \leq \dots \leq f(x_{\lambda:\lambda})$$



$$g(f(x_{1:\lambda})) \leq g(f(x_{2:\lambda})) \leq \dots \leq g(f(x_{\lambda:\lambda}))$$

for all $g : \text{Im}(f) \rightarrow \mathbb{R}$ strictly increasing

A Template for Comparison-based Stochastic Search

Define $\{P_\theta : \theta \in \Theta\}$, a family of probability distributions on \mathbb{R}^n

Generic template to optimize $f : \mathbb{R}^n \rightarrow \mathbb{R}$

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

While not terminate

1. Sample x_1, \dots, x_λ according to P_θ
2. Evaluate x_1, \dots, x_λ on f
3. Rank the solutions and find π the permutation such

$$f(x_{\pi(1)}) \leq f(x_{\pi(2)}) \leq \dots \leq f(x_{\pi(\lambda)})$$

4. Update parameters $\theta \leftarrow F(\theta, x_1, \dots, x_\lambda, \pi)$

How to update the different parameters m, σ, \mathbf{C} ?

1. Adapting the mean m
- 2. Adapting the step-size σ**
3. Adapting the covariance matrix \mathbf{C}

Why Step-size Adaptation?

Assume a $(1+1)$ -ES algorithm with fixed step-size σ (and $C = I_d$) optimizing the function $f(x) = \sum_{i=1}^n x_i^2 = \|x\|^2$.

Initialize \mathbf{m}, σ

While (stopping criterion not met)
sample new solution:

$$\mathbf{x} \leftarrow \mathbf{m} + \sigma \mathcal{N}(0, I_d)$$

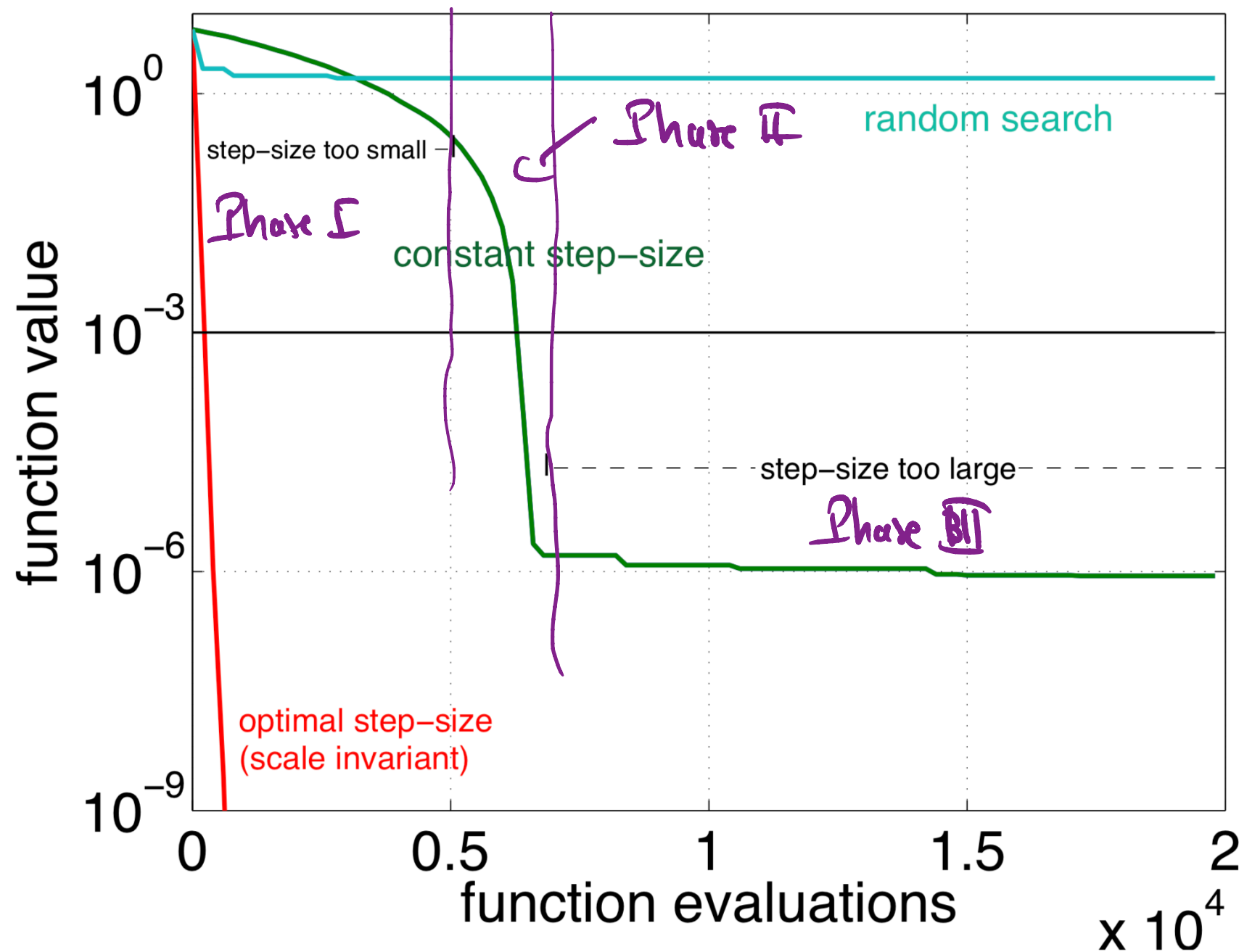
if $f(\mathbf{x}) \leq f(\mathbf{m})$

$$\mathbf{m} \leftarrow \mathbf{x}$$

$$\left(\begin{array}{l} \sigma \leftarrow \sigma \times 1,5 \\ \text{else} \quad \sigma \leftarrow \sigma \times (1,5)^{-1/4} \end{array} \right) \rightarrow \frac{1}{5} \text{ success-rule}$$

What will happen if you look at the convergence of $f(m)$?

Why Step-size Adaptation?



(1+1)-ES
(red & green)

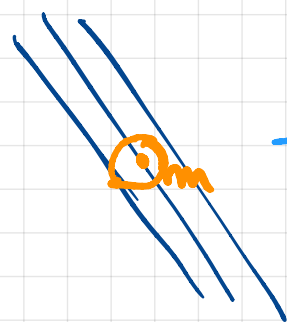
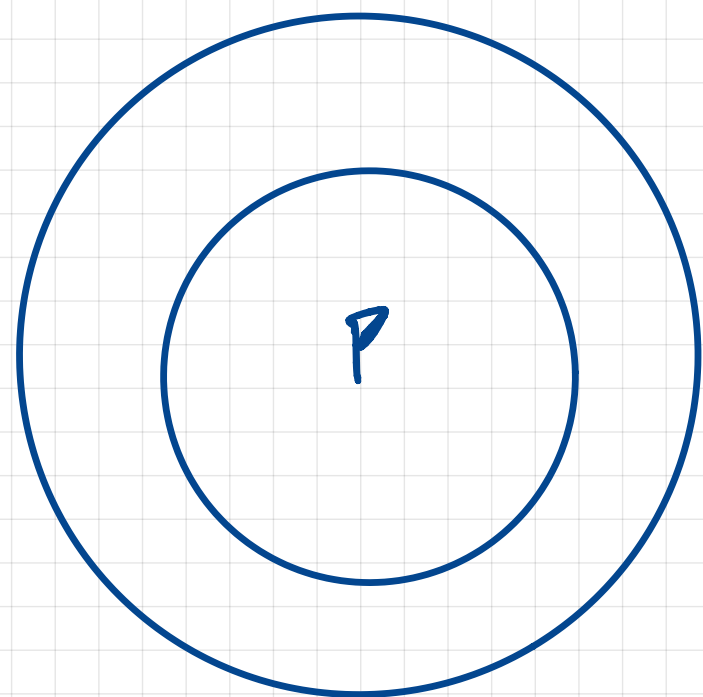
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

in $[-2.2, 0.8]^n$
for $n = 10$

red curve: (1+1)-ES with optimal step-size (see later)

green curve: (1+1)-ES with constant step-size ($\sigma = 10^{-3}$)

Phase I



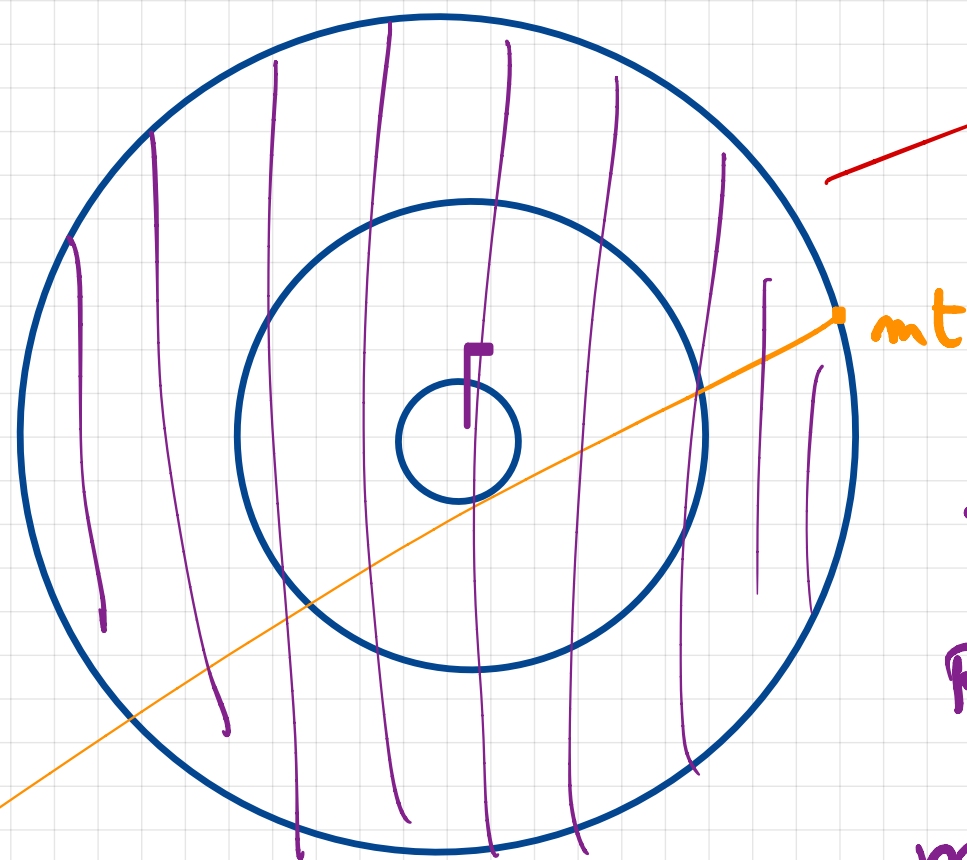
→ sample very close to
 $m_t = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ because step-size
very small

I progress but
slowly because
step-size is small.

\mathbb{P} (sample better solution)
 $\approx \frac{1}{2}$

↳ Step-size too small
with bigger, faster progress

Phase III



I would like to increase the step-size

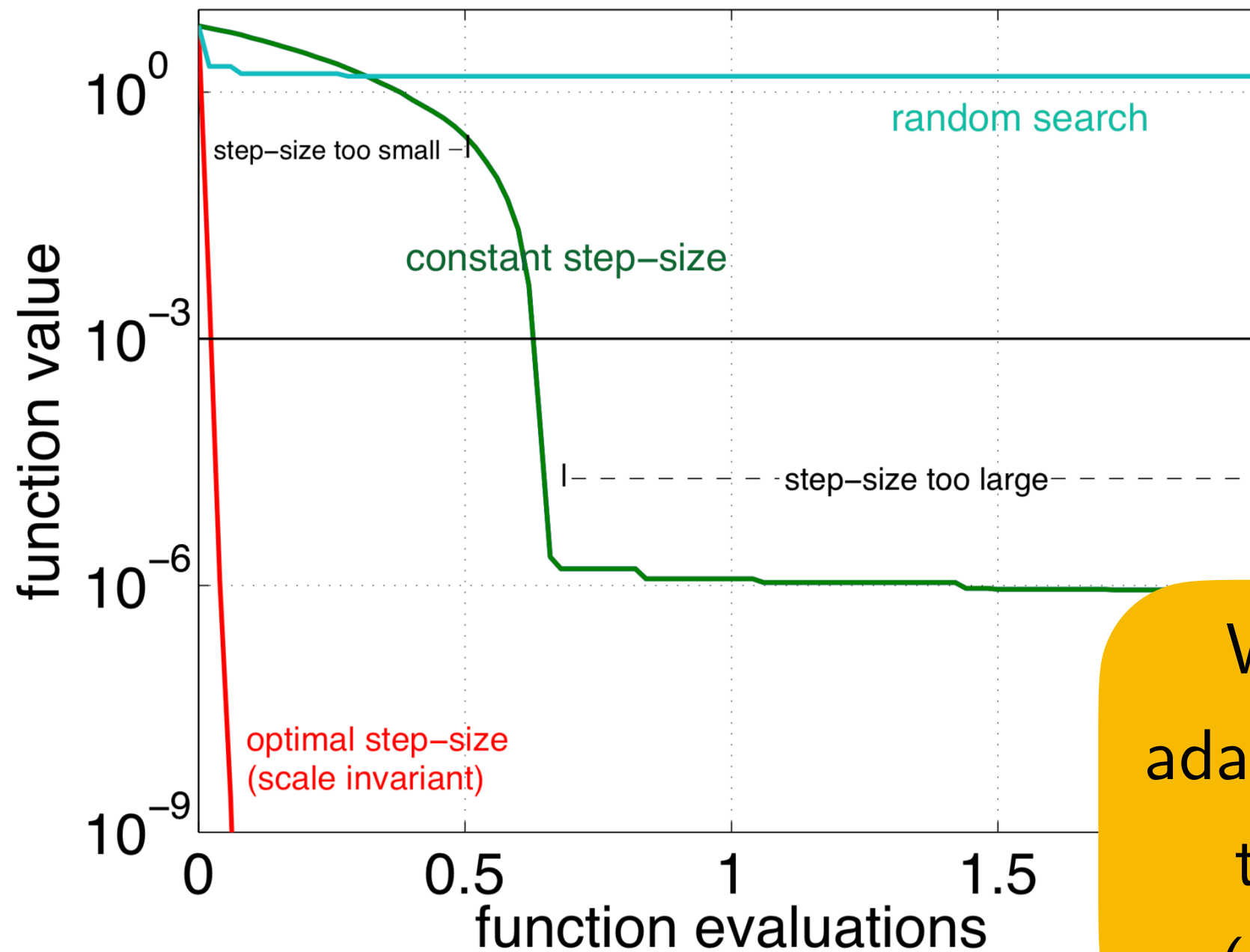
$$\{x \mid f(x) \leq f(mt)\}$$

Region where we need to sample to progress.

Probability to sample in this region is very small, because $\sigma \gg \|mt - x^*\|$

Phase II : The step size has the right order of magnitude compared to $\|mt - x\|$, progress is close to optimal.

Why Step-size Adaptation?



(1+1)-ES
(red & green)

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

We need step-size adaptation to approach the optimum fast (converge linearly)

red curve: (1+1)-ES with optimal step-size (see later)

green curve: (1+1)-ES with constant step-size ($\sigma = 10^{-3}$)

Methods for Step-size Adaptation

1/5th success rule, typically applied with “+” selection

[Rechenberg, 73][Schumer and Steiglitz, 78][Devroye, 72]

σ -self adaptation, applied with “,” selection

[Schwefel, 81]

random variation is applied to the step-size and the better one, according to the objective function value, is selected

path-length control or Cumulative step-size adaptation (CSA), applied with “,” selection

[Ostermeier et al. 84][Hansen, Ostermeier, 2001]

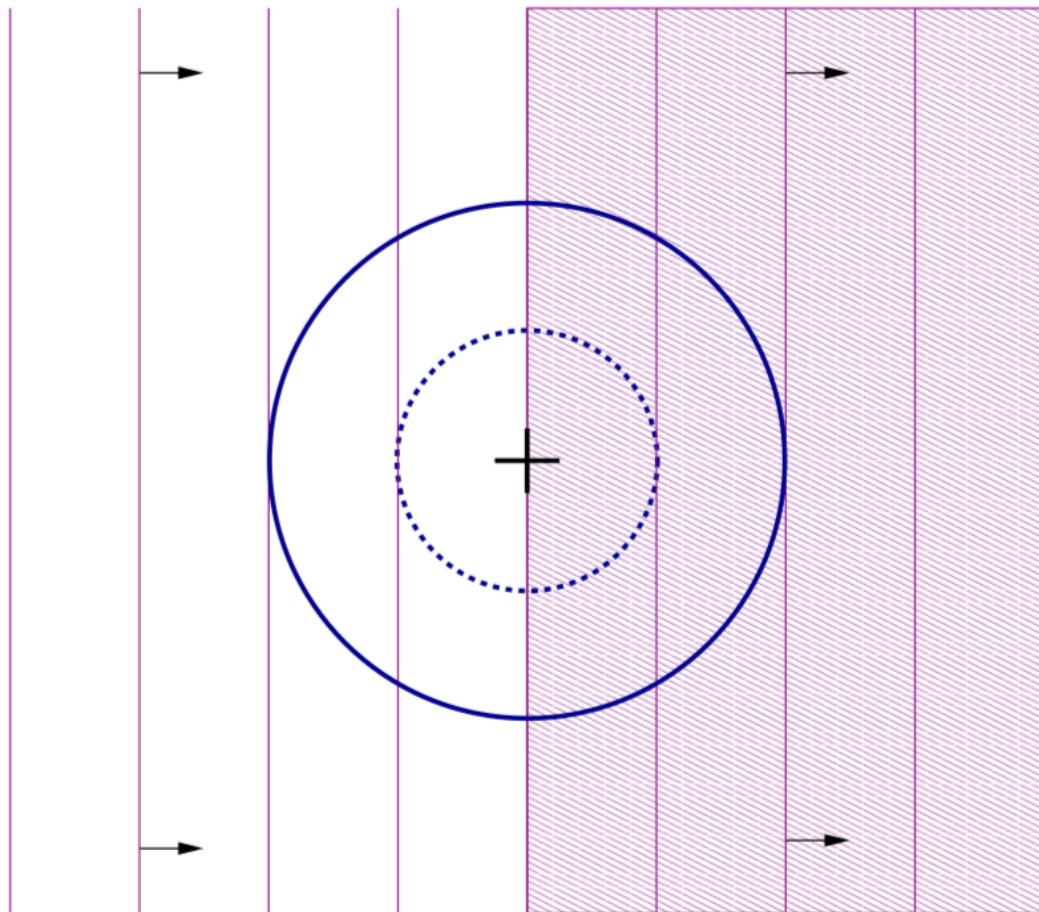
two-point adaptation (TPA), applied with “,” selection

[Hansen 2008]

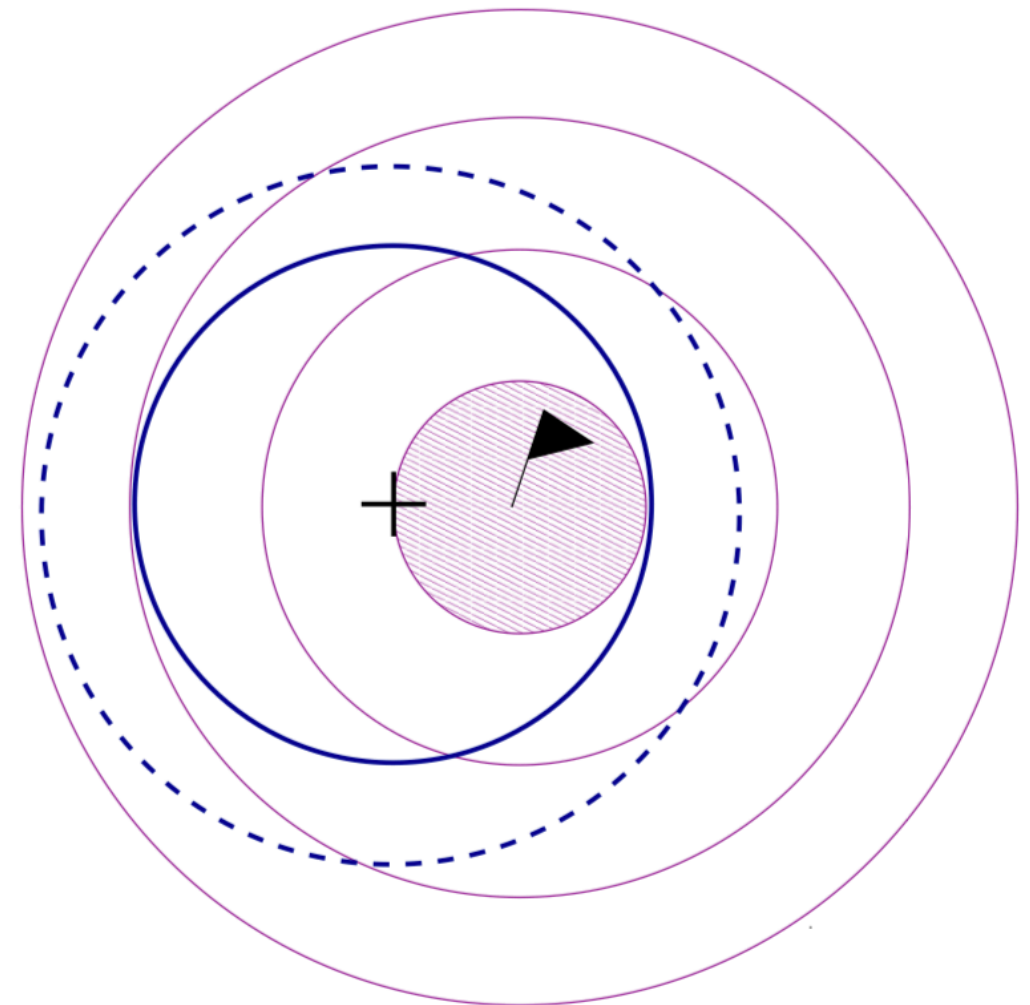
test two solutions in the direction of the mean shift, increase or decrease accordingly the step-size

Step-size control: 1/5th Success Rule

$$\mathbf{x} = (x_1, x_2)$$
$$f(\mathbf{x}) = -x_1$$

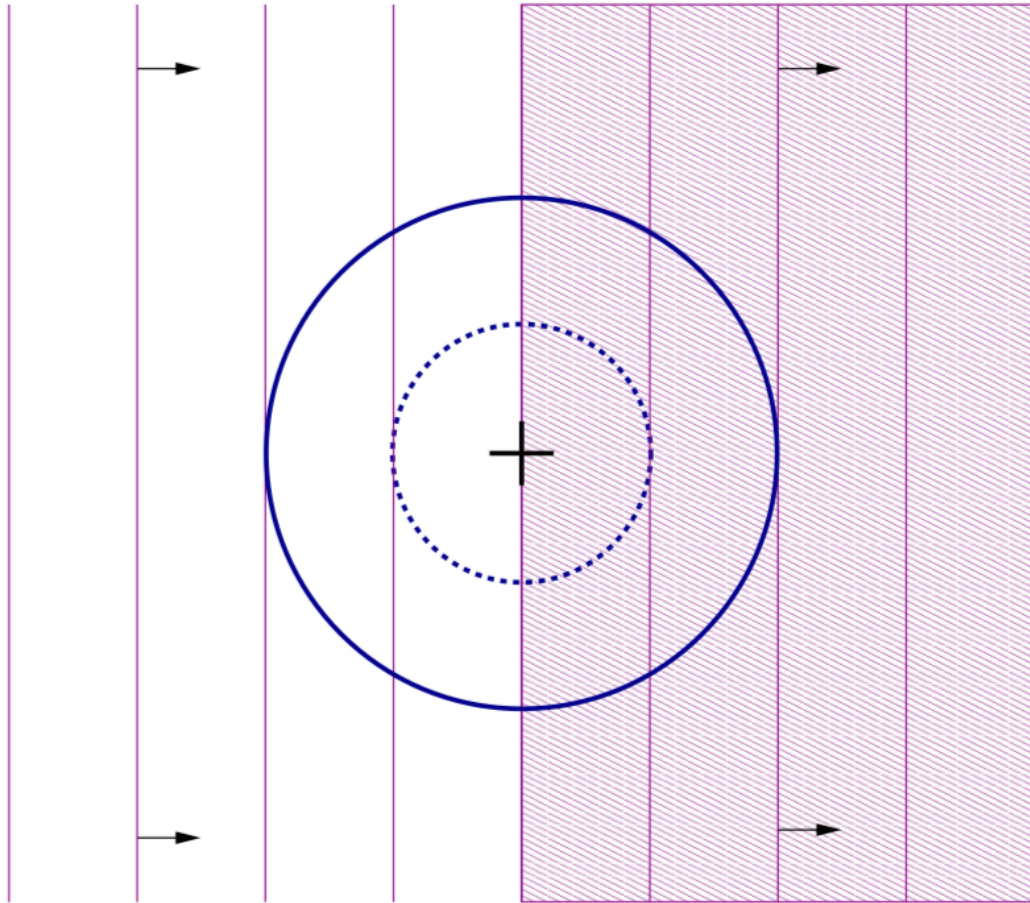


↓
increase σ



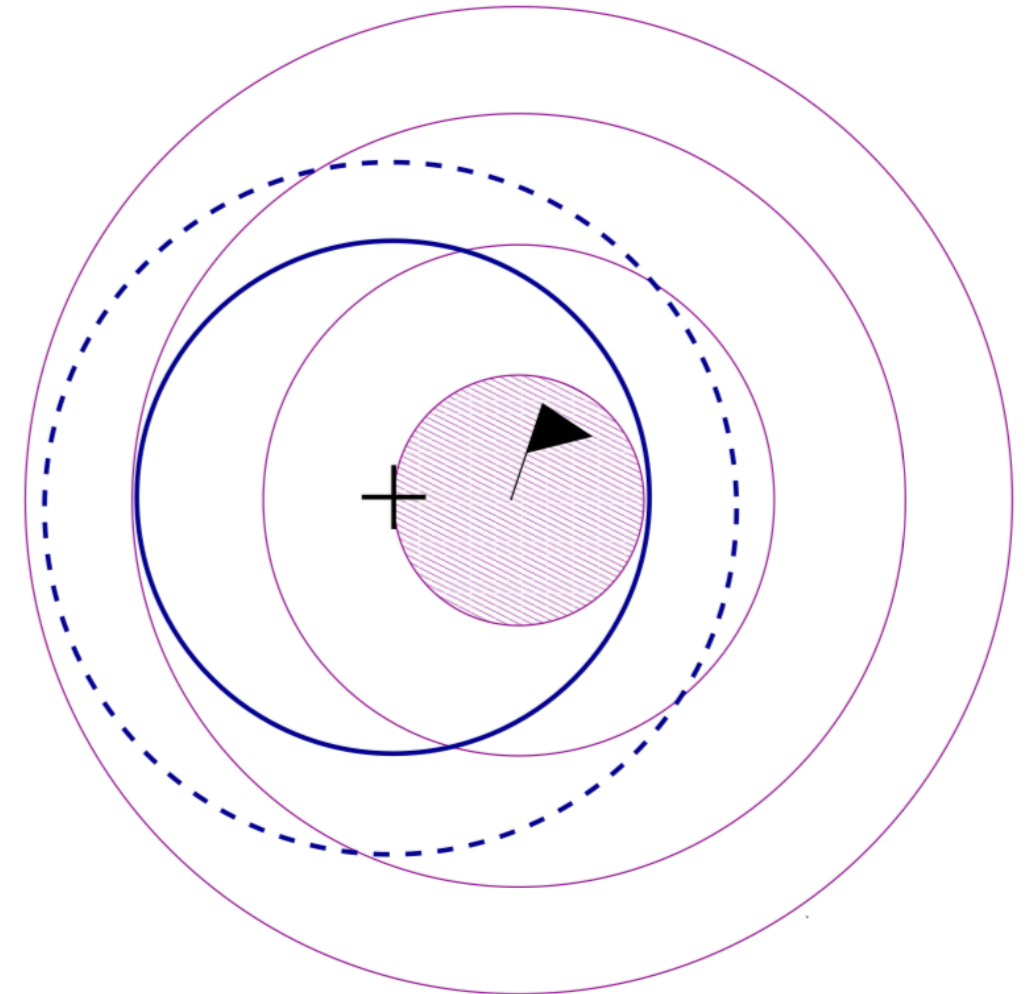
↓
decrease σ

Step-size control: 1/5th Success Rule



Probability of success (p_s)

$1/2$



Probability of success (p_s)

$1/5$

“too small”

Step-size control: 1/5th Success Rule

probability of success per iteration:

$$p_s = \frac{\text{\#candidate solutions better than } m}{\text{\#candidate solutions}}$$

$$\sigma \leftarrow \sigma \times \exp\left(\frac{1}{3} \times \frac{p_s - p_{\text{target}}}{1 - p_{\text{target}}}\right)$$

Increase σ if $p_s > p_{\text{target}}$
Decrease σ if $p_s < p_{\text{target}}$

(1 + 1)-ES

$$p_{\text{target}} = 1/5$$

IF *offspring better parent* [$f(\mathbf{x}) \leq f(\mathbf{m})$]

$$p_s = 1, \sigma \leftarrow \sigma \times \exp(1/3)$$

ELSE

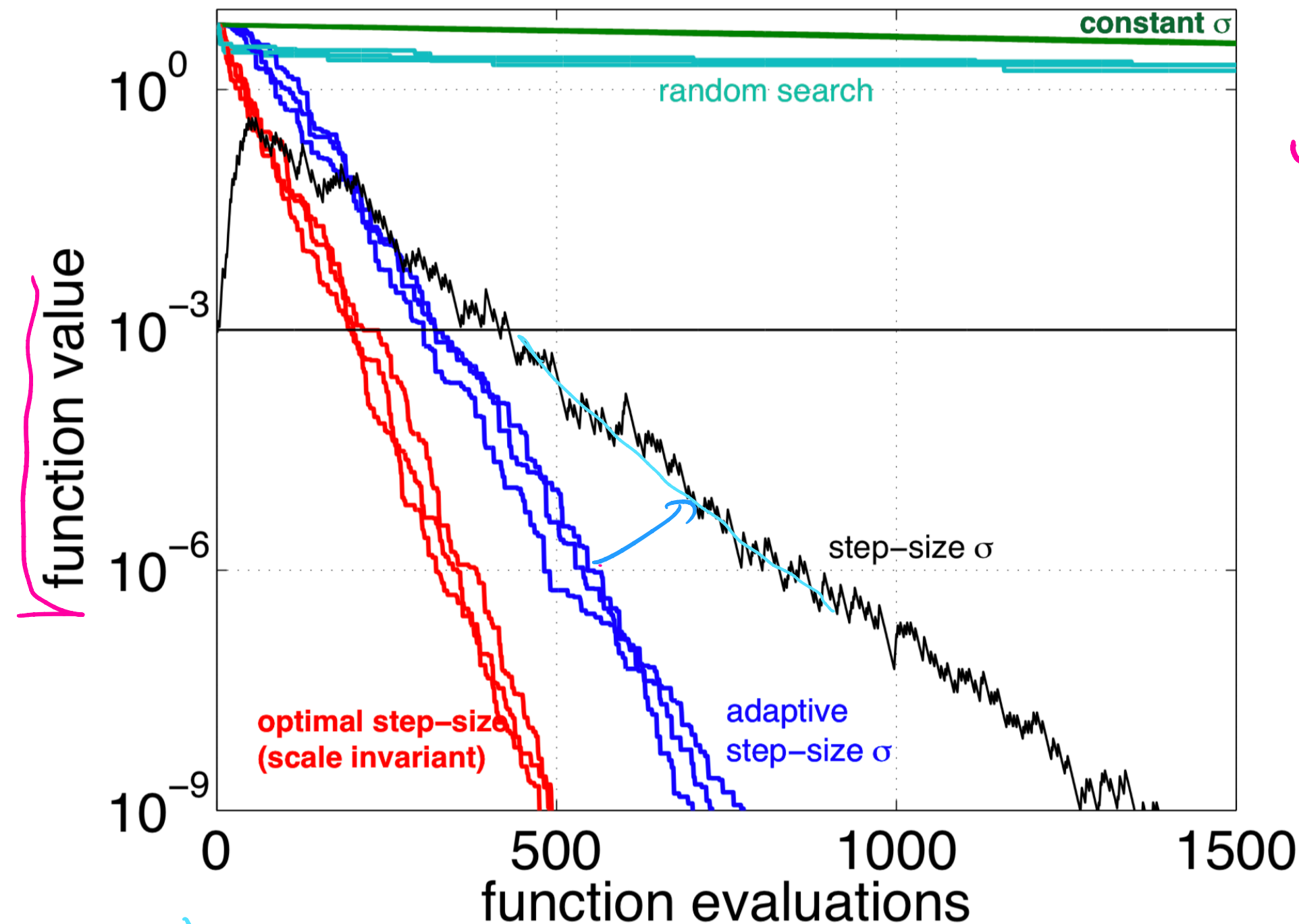
$$p_s = 0, \sigma \leftarrow \sigma / \exp(1/3)^{1/4}$$

1,5 [in exercise]

1,5^{1/4}

(1+1)-ES with One-fifth Success Rule - Convergence

(1 + 1)-ES with one-fifth success rule (blue)



If we display $\sqrt{f(x)} = \|x\|$ instead of $\|x\|^2$, then the blue convergence graph would have the same slope as the step-size slope (black).

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

in $[-0.2, 0.8]^n$
for $n = 10$

Linear convergence

$(mt, \sigma t)_{t \in \mathbb{N}}$
 $\log \|mt\|^2$
 $\log \sigma t$
 $\log \|mt\|^2$

The convergence observed on the previous slide is called linear convergence. We formalize it as

$\forall m_0$
 $\exists \delta_0$

$$\frac{1}{t} \ln \|m_t - x^*\| \xrightarrow{t \rightarrow \infty} -c$$

convergence rate

$$c > 0$$

(almost surely).

$-c$: slope of the black graph
slope of blue graph if we
display \sqrt{f}

Path Length Control - Cumulative Step-size Adaptation (CSA)

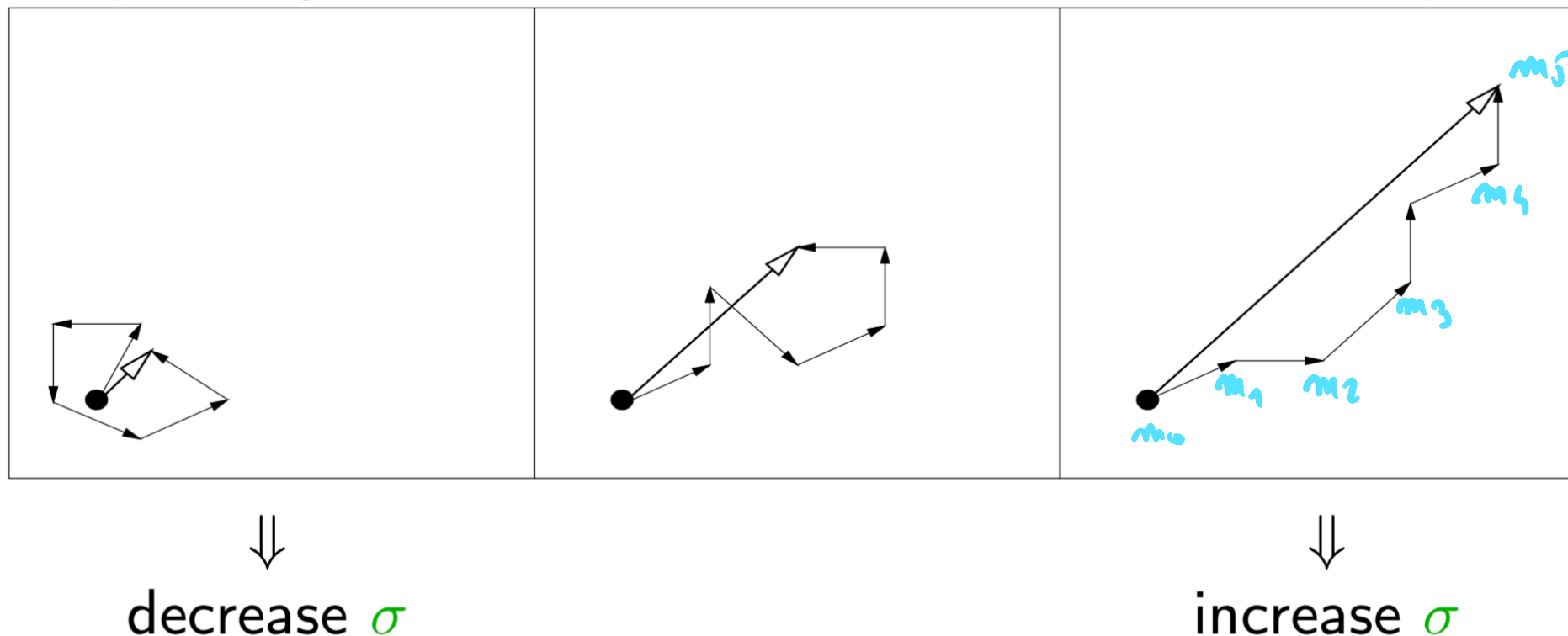
step-size adaptation used in the $(\mu/\mu_w, \lambda)$ -ES algorithm framework (in CMA-ES in particular)

Main Idea:

$$\begin{aligned} \mathbf{x}_i &= \mathbf{m} + \sigma \mathbf{y}_i \\ \mathbf{m} &\leftarrow \mathbf{m} + \sigma \mathbf{y}_w \end{aligned}$$

Measure the length of the *evolution path*

the pathway of the mean vector \mathbf{m} in the iteration sequence



Sampling of solutions, notations as on slide “The $(\mu/\mu, \lambda)$ -ES - Update of the mean vector” with \mathbf{C} equal to the identity.

Initialize $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, evolution path $\mathbf{p}_\sigma \in \mathbb{R}^n = \mathbf{0}$,
set $c_\sigma \approx 4/n$, $d_\sigma \approx 1$.

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \mathbf{y}_w \quad \text{where } \mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \quad \text{update mean}$$

$$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \underbrace{\sqrt{1 - (1 - c_\sigma)^2}}_{\text{accounts for } 1 - c_\sigma} \underbrace{\sqrt{\mu_w}}_{\text{accounts for } w_i} \mathbf{y}_w$$

$$\sigma \leftarrow \sigma \times \underbrace{\exp \left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E} \|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)}_{>1 \iff \|\mathbf{p}_\sigma\| \text{ is greater than its expectation}} \quad \text{update step-size}$$

Handwritten notes:

$$x_i = \mathbf{m} + \sigma \mathbf{y}_i'$$

$$f(x_{1:\lambda}) \leq \dots \leq f(x_{\lambda:\lambda})$$

$$x_{i:\lambda} = \mathbf{m} + \sigma \mathbf{y}_{i:\lambda}$$

$\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

The CSA update with time index notation writes:

$$m_{t+1} = m_t + \sigma_t \sum_{i=1}^{\mu} w_i y_{t+1}^{i:\lambda}$$

$$p_{t+1}^{\sigma} = (1-\alpha) p_t^{\sigma} + \sqrt{1-(1-\alpha)^2} \sqrt{\mu w} x_{t+1}^w$$

$$\sigma_{t+1} = \sigma_t \exp \left(\frac{\alpha}{d\sigma} \left(\frac{\|p_{t+1}^{\sigma}\|}{\mathbb{E}(\|N(0, Id)\|)} - 1 \right) \right)$$

$$x_i = m_t + \sigma_t x_{t+1}^i$$

$$\{y_{t+1}^i, i=1, \dots, \lambda\}$$

$$\sum_{i=1}^{\mu} w_i y_{t+1}^{i:\lambda} \stackrel{iid}{\sim} N(0, Id)$$

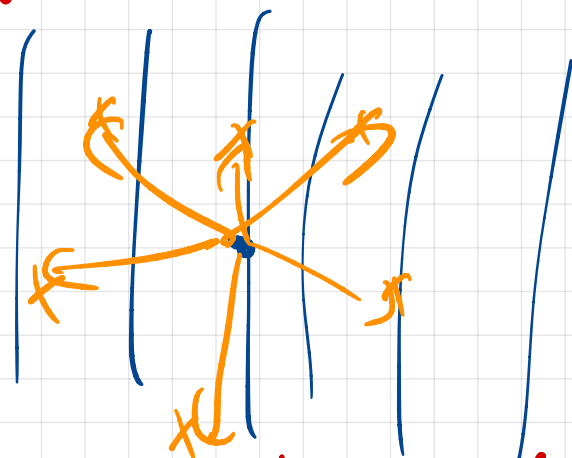
$$\mu w = \frac{1}{\sum w_i^2}$$

When $f(x) = \text{random}$ [each time we evaluate a point on f , it returns a random number independant from the previous ones and identically distributed]

then $\{y_{t+1}^{i:\lambda}, i=1, \dots, \lambda\}$ are iid $\sim N(0, Id)$

$$y_{t+1}^i \sim N(0, Id)$$

is not true on other functions. The selection



on f , generally changes the distribut°

If $y_{t+1}^{i:\lambda} \sim \mathcal{N}(0, Id)$, then $\sqrt{\mu w} y_w^{t+1} = \frac{1}{\sqrt{\sum w_i^2}} \sum_{i=1}^{\mu} w_i y_{t+1}^{i:\lambda} \sim \mathcal{N}(0, Id)$
(Hone Proof)

under random selection

If $p_0^t \sim \mathcal{N}(0, Id)$, \checkmark since $\sqrt{\mu w} y_w^{t+1} \sim \mathcal{N}(0, Id)$

then

$$p_0^{t+1} = (1 - c\sigma) p_0^t + \sqrt{1 - (1 - c\sigma)^2} \sqrt{\mu w} y_w^{t+1} \sim \mathcal{N}(0, Id)$$

(Hone Proof)

Then : under random selection :

$$\sigma_{t+1} = \sigma_t \exp \left(\frac{c\sigma}{d\sigma} \left(\frac{\|p_{t+1}^\sigma\|}{\mathbb{E}(\|\mathcal{N}(0, Id)\|)} - 1 \right) \right)$$

$$\log \sigma_{t+1} = \log \sigma_t + \frac{c\sigma}{d\sigma} \left(\frac{\|p_{t+1}^\sigma\|}{\mathbb{E} \|N(0, Id)\|} - 1 \right)$$

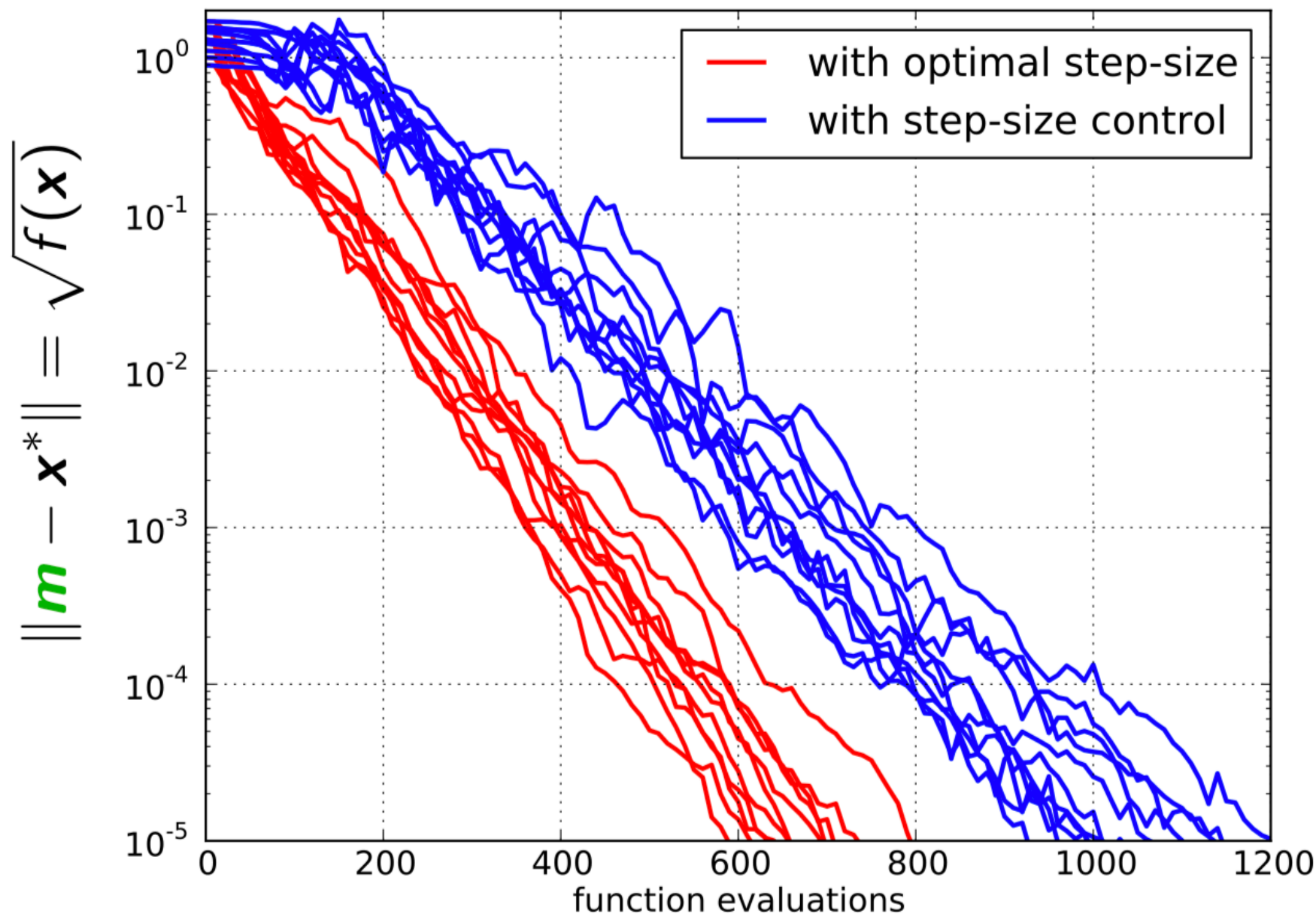
$$\mathbb{E} [\log \sigma_{t+1} \mid m_t, \sigma_t] = \log \sigma_t + \underbrace{\frac{c\sigma}{d\sigma} \left(\frac{\mathbb{E}(\|p_{t+1}^\sigma\|)}{\mathbb{E}(\|N(0, Id)\|)} - 1 \right)}_{=0}$$

$$\mathbb{E} [\log(\sigma_{t+1}) \mid m_t, \sigma_t] = \log \sigma_t$$

↳ Formal explanation of log step-size constant under random selection.

Convergence of $(\mu/\mu_w, \lambda)$ -CSA-ES

2x11 runs



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

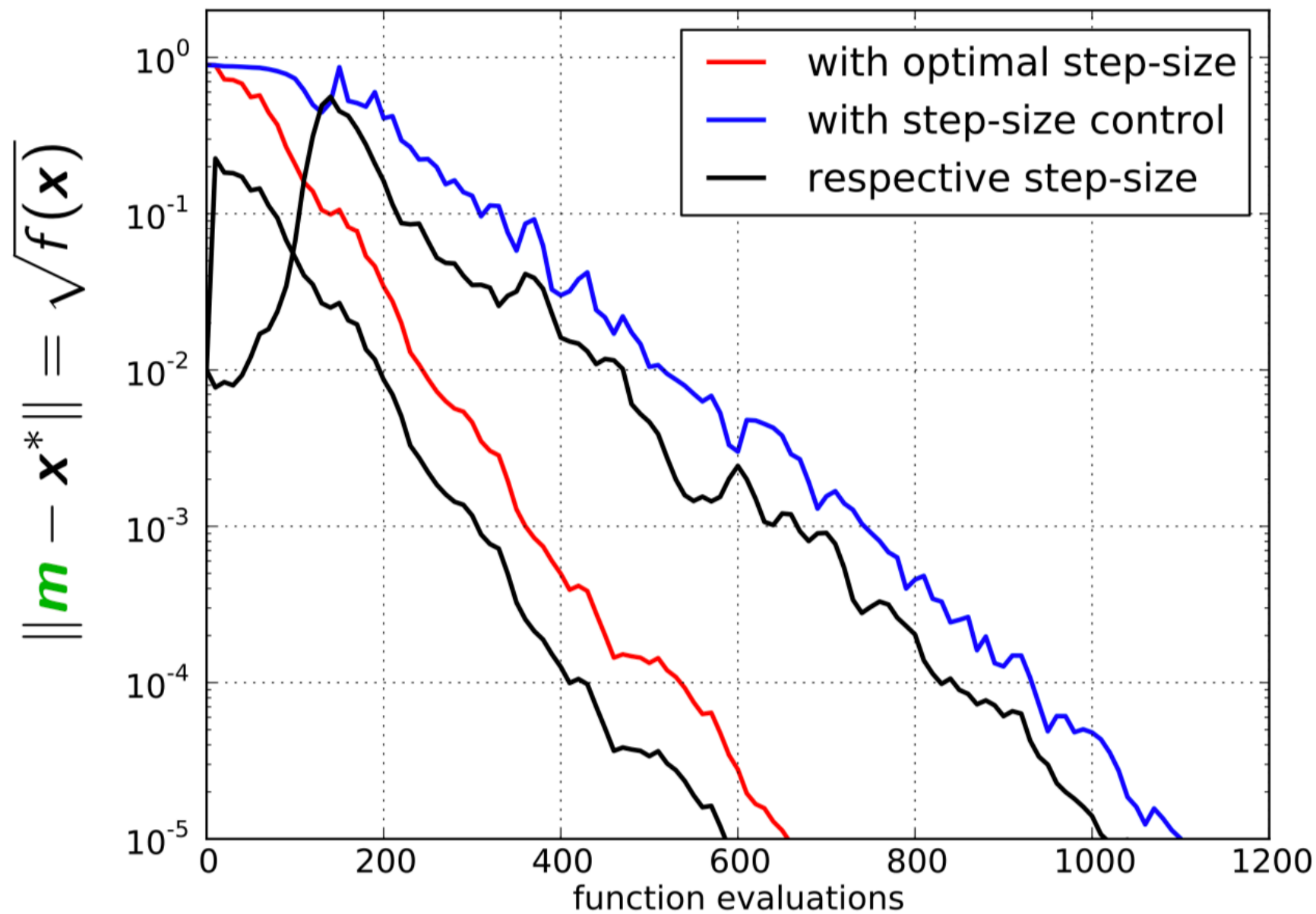
for $n = 10$

and

$$\mathbf{x}^0 \in [-0.2, 0.8]^n$$

with optimal versus adaptive step-size σ with too small initial σ

Convergence of $(\mu/\mu_w, \lambda)$ -CSA-ES



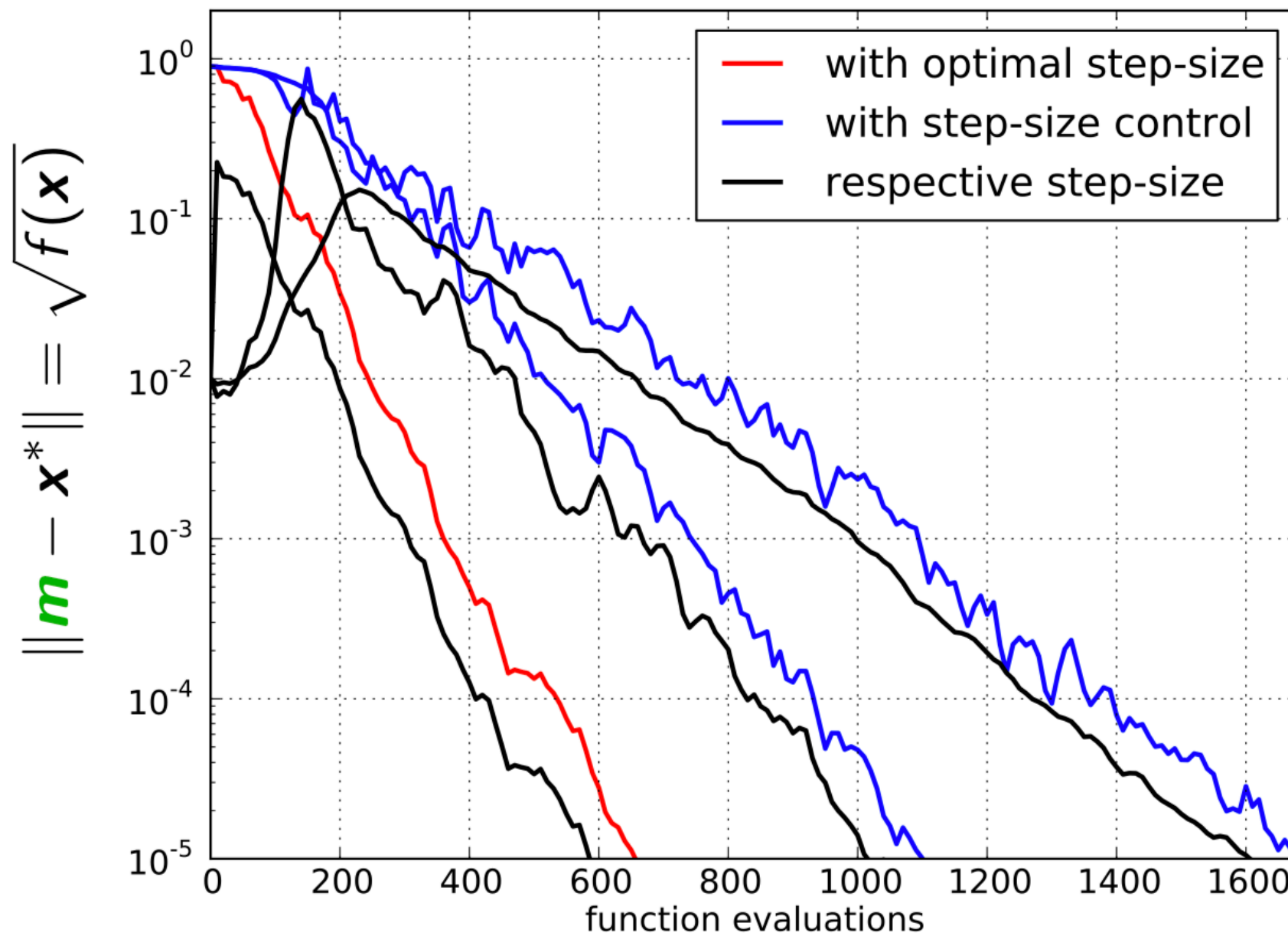
$$f(\mathbf{x}) = \sum_{i=1}^n x_i^2$$

for $n = 10$
and
 $\mathbf{x}^0 \in [-0.2, 0.8]^n$

comparing number of f -evals to reach $\|\mathbf{m}\| = 10^{-5}$: $\frac{1100-100}{650} \approx 1.5$

Note: initial step-size taken too small ($\sigma_0 = 10^{-2}$) to illustrate the step-size adaptation

Convergence of $(\mu/\mu_w, \lambda)$ -CSA-ES



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

for $n = 10$

and

$$\mathbf{x}^0 \in [-0.2, 0.8]^n$$

comparing optimal versus default damping parameter d_σ :

$$\frac{1700}{1100} \approx 1.5$$

Optimal Step-size - Lower-bound for Convergence Rates

In the previous slides we have displayed some runs with “optimal” step-size.

Optimal step-size relates to step-size proportional to the distance to the optimum: $\sigma_t = \sigma \|x - x^\star\|$ where x^\star is the optimum of the optimized function (with σ properly chosen).

The associated algorithm is not a real algorithm (as it needs to know the distance to the optimum) but it gives bounds on convergence rates and allows to compute many important quantities.

The goal for a step-size adaptive algorithm is to achieve convergence rates close to the one with optimal step-size

We will formalize this in the context of the $(1+1)$ -ES. Similar results can be obtained for other algorithm frameworks.

Optimal Step-size - Bound on Convergence Rate - (1+1)-ES

Consider a (1+1)-ES algorithm with **any step-size adaptation** mechanism:

$$X_{t+1} = \begin{cases} X_t + \sigma_t \mathcal{N}_{t+1} & \text{if } f(X_t + \sigma_t \mathcal{N}_{t+1}) \leq f(X_t) \\ X_t & \text{otherwise} \end{cases}$$

with $\{\mathcal{N}_t, t \geq 1\}$ i.i.d. $\sim \mathcal{N}(0, I_d)$

equivalent writing:

$$X_{t+1} = X_t + \sigma_t \mathcal{N}_{t+1} \mathbf{1}_{\{f(X_t + \sigma_t \mathcal{N}_{t+1}) \leq f(X_t)\}}$$

Bound on Convergence Rate - (1+1)-ES

Theorem: For any objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, for any $y^\star \in \mathbb{R}^n$

$$E[\|X_{t+1} - y^\star\|] \geq E[\|X_t - y^\star\|] - \tau \quad \text{lower bound}$$

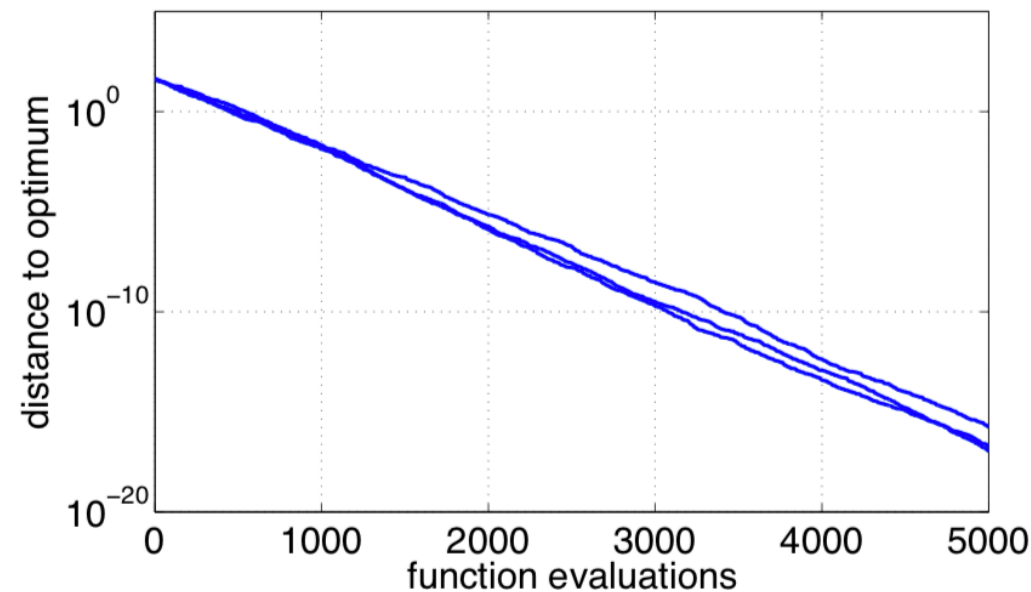
where $\tau = \max_{\sigma \in \mathbb{R}^+} \underbrace{E[\ln^- \|e_1 + \sigma \mathcal{N}\|]}_{=: \varphi(\sigma)}$ with $e_1 = (1, 0, \dots, 0)$

Theorem: The convergence rate lower-bound is reached on spherical functions $f(x) = g(\|x - x^\star\|)$ (with $g : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ strictly increasing) and step-size proportional to the distance to the optimum $\sigma_t = \sigma_{\text{opt}} \|x - x^\star\|$ with σ_{opt} such that $\varphi(\sigma_{\text{opt}}) = \tau$.

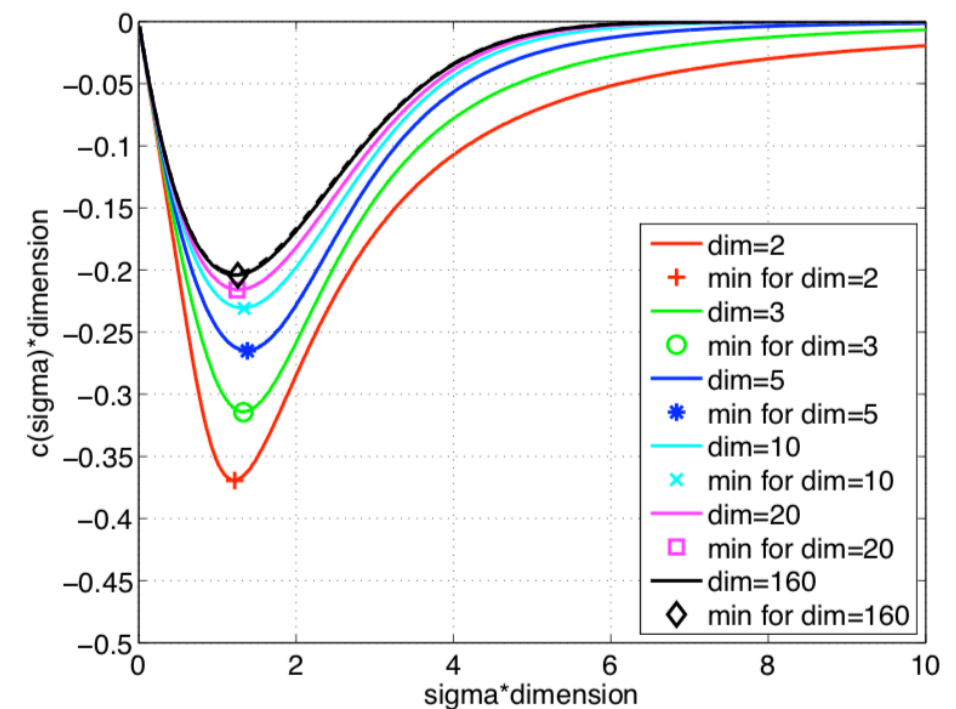
Log-Linear Convergence of scale-invariance step-size ES

Theorem: The (1+1)-ES with step-size proportional to the distance to the optimum $\sigma_t = \sigma \|x\|$ converges (log)-linearly on the sphere function $f(x) = g(\|x\|)$ almost surely:

$$\frac{1}{t} \ln \frac{\|X_t\|}{\|X_0\|} \xrightarrow{t \rightarrow \infty} -\varphi(\sigma) =: \text{CR}_{(1+1)}(\sigma)$$



$$n = 20 \text{ and } \sigma = 0.6/n$$



Asymptotic Results ($n \rightarrow \infty$)

Theorem

Let $\sigma > 0$, the convergence rate of the (1+1)-ES with scale-invariant step-size on spherical functions satisfies at the limit

$$\lim_{n \rightarrow \infty} n \times \text{CR}_{(1+1)} \left(\frac{\sigma}{n} \right) = \frac{-\sigma}{\sqrt{2\pi}} \exp \left(-\frac{\sigma^2}{8} \right) + \frac{\sigma^2}{2} \Phi \left(-\frac{\sigma}{2} \right)$$

where Φ is the cumulative distribution of a normal distribution.

optimal convergence rate decreases to zero like $\frac{1}{n}$

