

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]

→ finite differences

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]

→ Exploits separability

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] → originally for discrete,

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 (typically ind).
(Handwritten: $x_1, \dots, x_\lambda \in \mathbb{R}^n$)
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?

Assume $n=1$, we can sample with a normal distribution.

Definition: A random variable $X: (\Omega, \mathcal{F}) \rightarrow \mathbb{R}$ is a normal distribution with mean m and variance σ^2 if its probability density function equals:

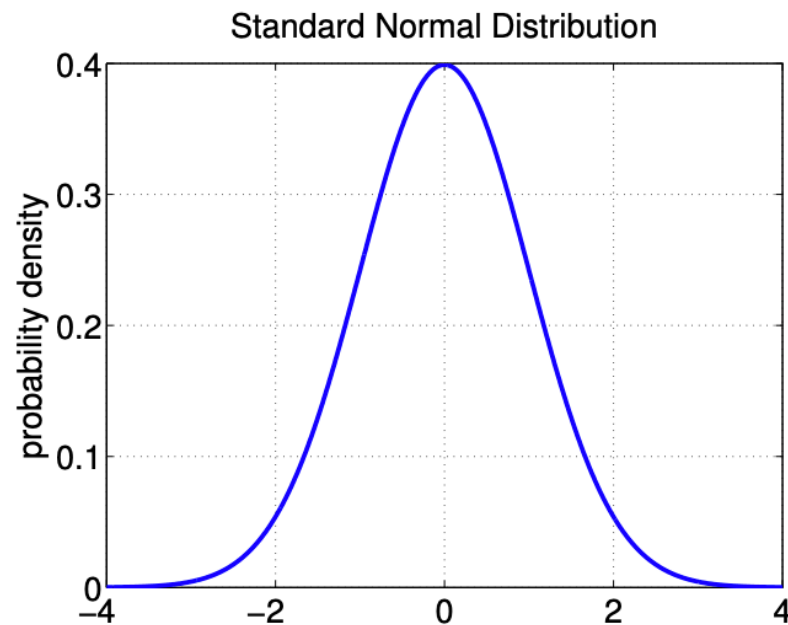
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right)$$

We denote $X \sim \mathcal{N}(m, \sigma^2)$

If $m=0$ and $\sigma=1$, we talk about standard normal distribution.

The following holds: $m + \sigma \mathcal{N}(0, 1) \sim \mathcal{N}(m, \sigma^2)$

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

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

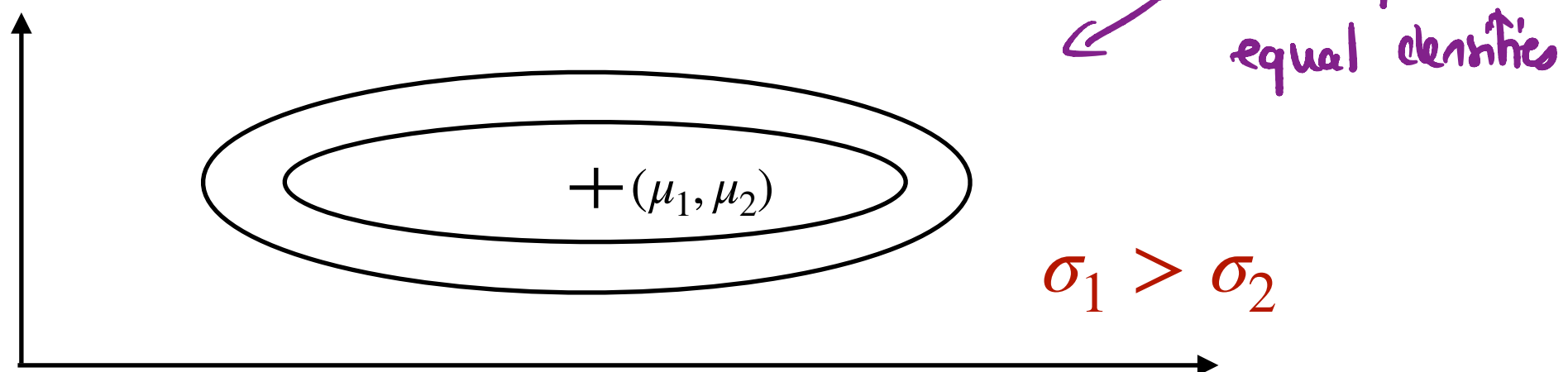
Assume $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ denote its density $p(x_1) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2\right)$

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

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

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 $a_1X_1 + \dots + a_nX_n = \langle a, X \rangle$ has a **normal distribution**.

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

In the n-dimensional case a multivariate normal distribution is **determined** by its mean vector **m** and covariance matrix **C**.

Reminder: Covariance matrix

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

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

where $\mu_i = E[X_i]$. Considering the operator that take the expectation of each entry of a matrix: $\Sigma = E[(X - \mu)^\top (X - \mu)]$.

$$\Sigma_{ii} = \text{Var}(X_i) = E(X_i^2) - E(X_i)^2$$

positive semidefinite

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

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

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

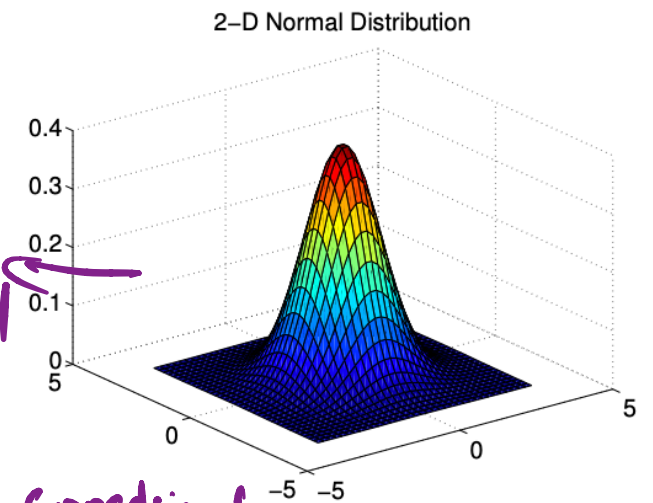
The **mean vector** m :

determines the displacement

is the value with the largest density

the distribution is symmetric around the mean

$C = I_d$
Line of equal
densities
= circle.
↳ both coordinates
are i.i.d.



$$\mathcal{N}(m, C) \sim m + \mathcal{N}(0, C) \sim m + C^{1/2} \mathcal{N}(0, I_d)$$

The **covariance matrix**:

determines the geometrical shape (see next slides)

$$C^{1/2} \stackrel{\text{st}}{=} C^{1/2} C^{1/2} = C$$

If C symmetric definite positive

$$C = U^T D U \quad \text{where } U^T U = \text{Id}$$

I can define $C^{1/2} = U^T \sqrt{D} U \rightarrow$ This is the unique sym square root.

D is diagonal with $D_{ii} > 0$ since C is definite positive.

$$\sqrt{D} = \begin{pmatrix} \sqrt{D_{11}} & & (0) \\ & \ddots & \\ (0) & & \sqrt{D_{nn}} \end{pmatrix}$$

Geometry of a Gaussian Vector

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

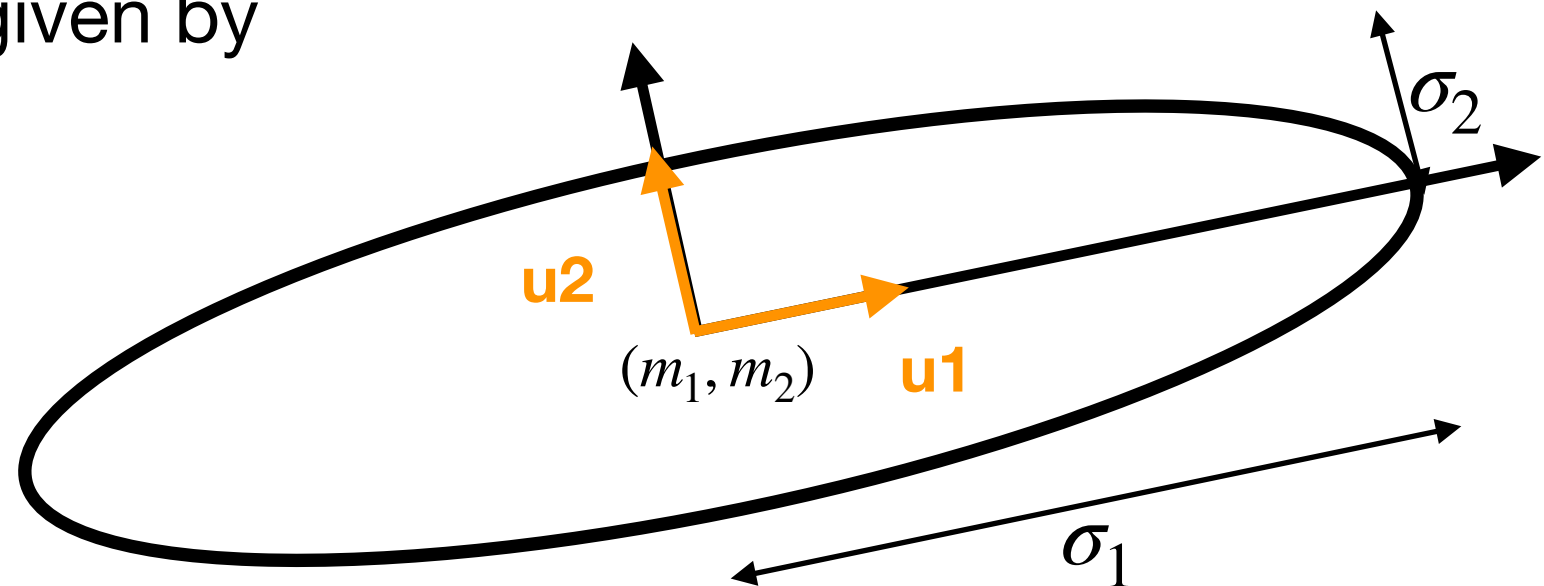
$$\{x \mid \Delta^2 = (x - m)^T C^{-1} (x - m) = \text{cst}\}$$

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

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

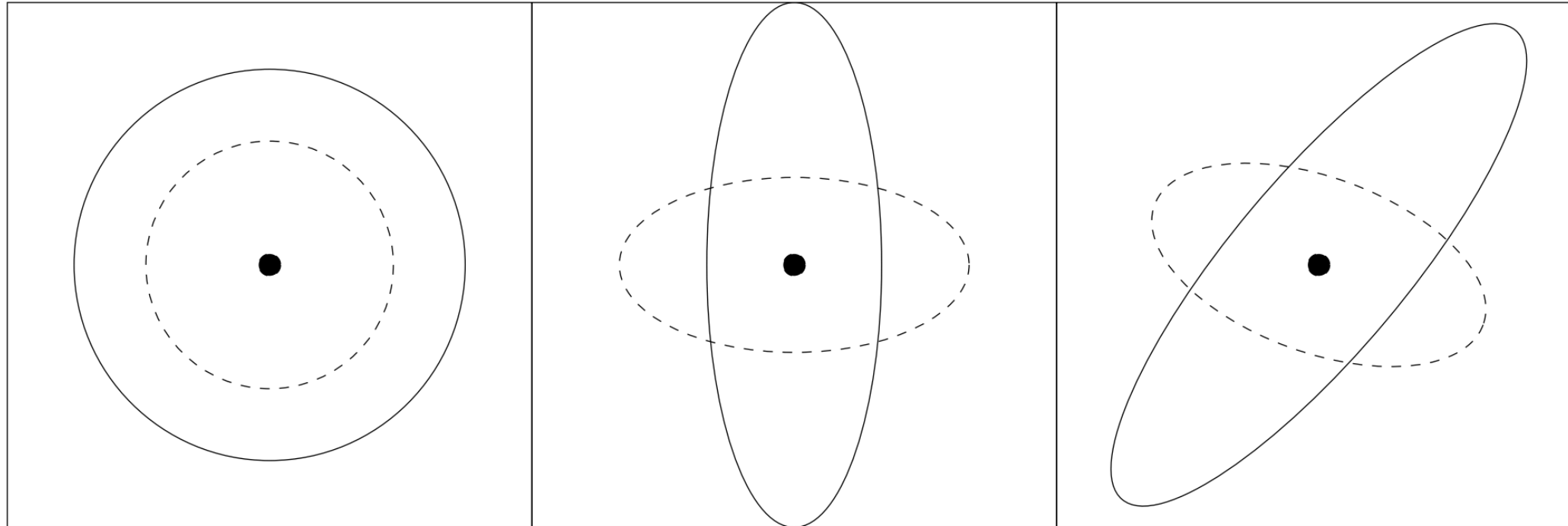
Let $Y = U^T(x - m)$, then in the coordinate system, (u_1, u_2) , the lines of equal densities are given by

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



... 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

Introduced in the 70's

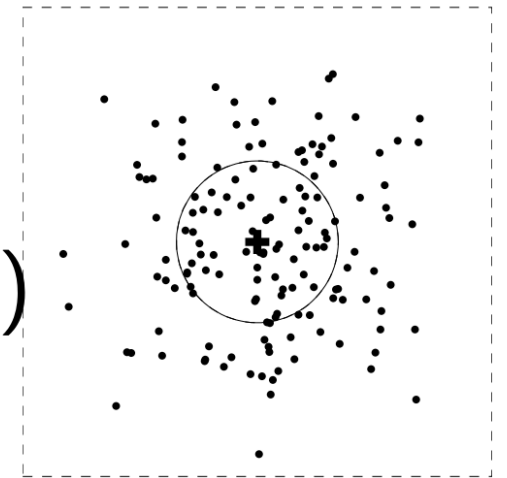
New search points are sampled normally distributed

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

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

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

- ▶ 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

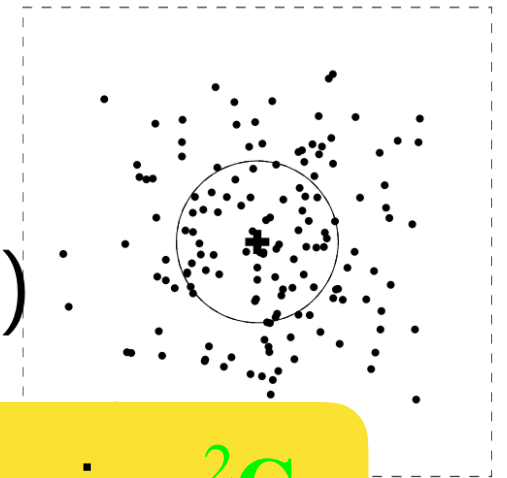
(2001, Hansen
& Ostermeier)

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}

How to adapt m ?

Assume $C = Id$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

5th
 x_{t+1}^1
 x^1

4th
 x_{t+1}^4
 x^4

m_t

1st
 x_{t+1}^2
 x^2

3rd
 x_{t+1}^3
 x^3

2nd
 x_{t+1}^5
 x^5

$$m_{t+1} = \frac{1}{Z} \sum \frac{1}{f(x_t^i)} x_t^i$$

normalise

problem in zero
→ Take exponential

$\exp(-\frac{1}{T} f(x_{t+1}^i))$
Temperature

Introduce a notation: $f(x_{t+1}^{1:\lambda}) \leq f(x_{t+1}^{2:\lambda}) \leq \dots \leq f(x_{t+1}^{n:\lambda})$

Other idea:

$$m_{t+1} = x_{t+1}^{1:\lambda}$$

→ $(1, \lambda) - EI$.

Take more points to have more information like half of them

Simpler setting:

Assume that at each iteration $\lambda = 1$

\hat{x}_{t+1}^1 • m_t

$$m_t = \begin{cases} \hat{x}_{t+1}^1 & \text{if } f(\hat{x}_{t+1}^1) \leq f(m_t) \\ m_t & \text{otherwise} \end{cases}$$

$$\rightarrow (1 + \epsilon) - \epsilon$$

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 $\mathbf{x}_i = \mathbf{m} + \sigma \underbrace{\mathbf{y}_i}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{C})}$

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

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

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

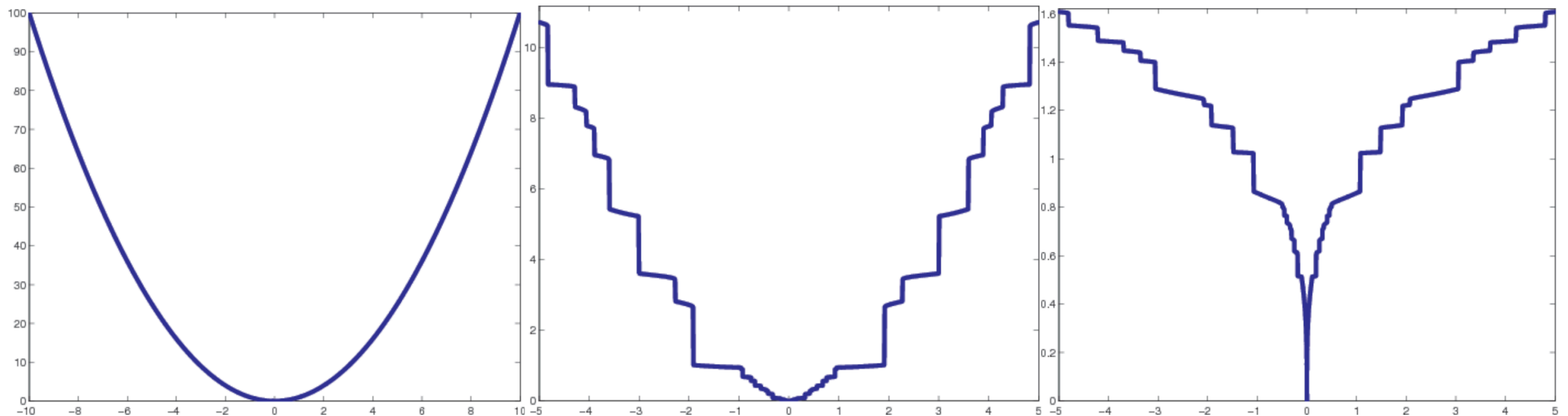
What changes in the previous slide if instead of optimizing f , we optimize $g \circ f$ where $g : \text{Im}(f) \rightarrow \mathbb{R}$ is strictly increasing?

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 $\theta \leftarrow F(\theta, x_1, \dots, x_\lambda, \pi)$

$\pi(1) : 1 : \lambda$
 \vdots
 $\pi(\lambda) : (\lambda : \lambda)$

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}

Exercise : Adaptive step-size algorithms

III Adaptive step-size algorithms

We are going to test the convergence of several algorithms on some test functions, in particular on the so-called sphere function

$$f_{\text{sphere}}(\mathbf{x}) = \sum_{i=1}^n \mathbf{x}_i^2 = \|\mathbf{x}\|^2$$

and the ellipsoid function

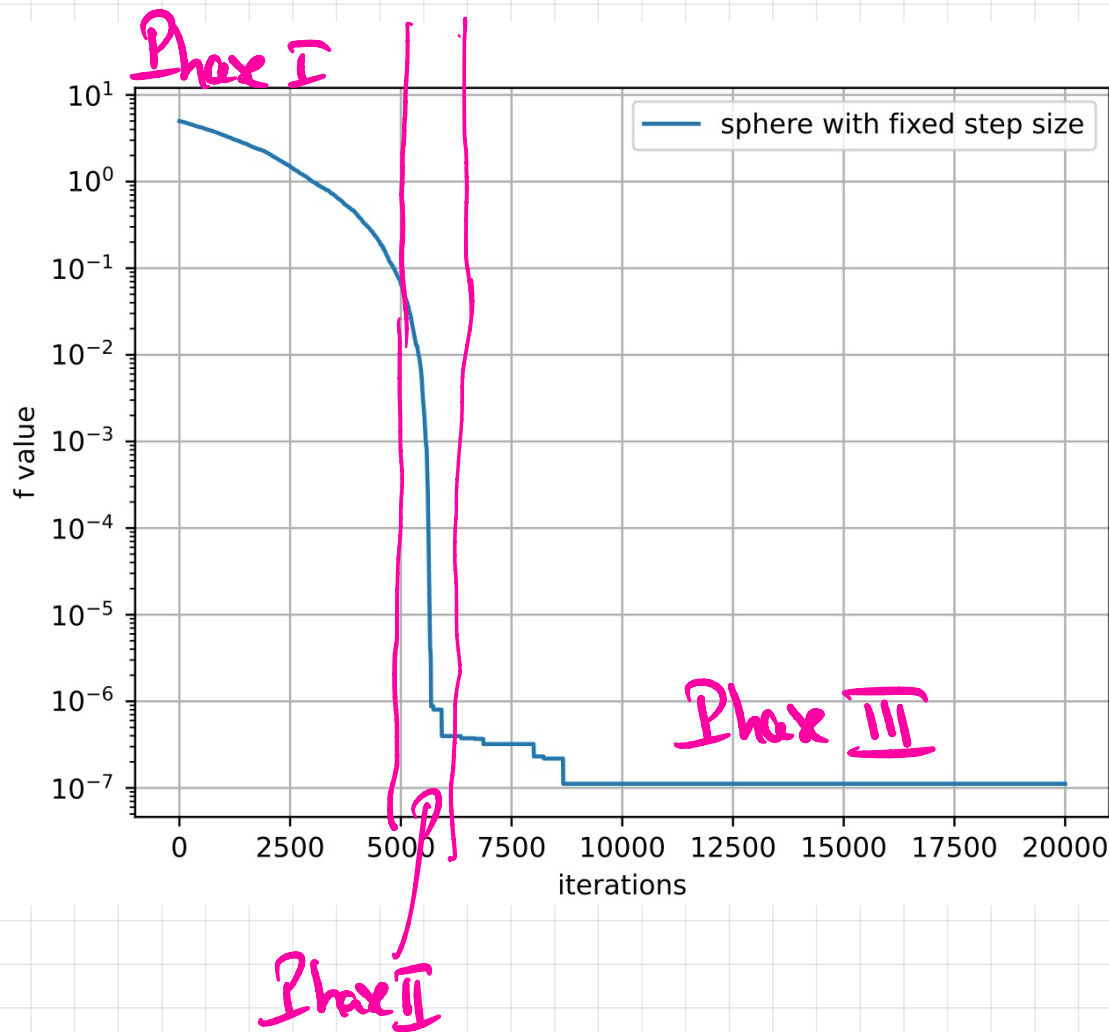
$$f_{\text{elli}}(\mathbf{x}) = \sum_{i=1}^n (100^{\frac{i-1}{n-1}} \mathbf{x}_i)^2 .$$

1. What is the condition number associated to the Hessian matrix of the functions above? Are the functions ill-conditioned?

f_{sphere} : Hessian = 2 Id \rightarrow cond = 1

f_{elli} : Hessian = $2 \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 10^4 \end{pmatrix}$ cond = $\frac{2 \cdot 10^4}{2 \cdot 1} = 10^4$

Convergence (1+1)-ES with fixed step-size.



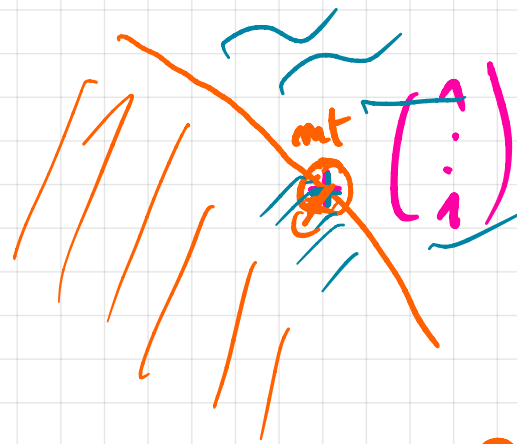
on $f_{\text{sphere}}(x) = \|x\|^2$

$$n = 5$$

$$x_0 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\sigma_0 = 10^{-3}$$

Phase II: step-size well adapted compared to $\| \text{mut} \|$



Phase I:

σ_t very small compared
to mt

→ Progress slowly

Probability of improvement

$$\approx \frac{1}{2}$$

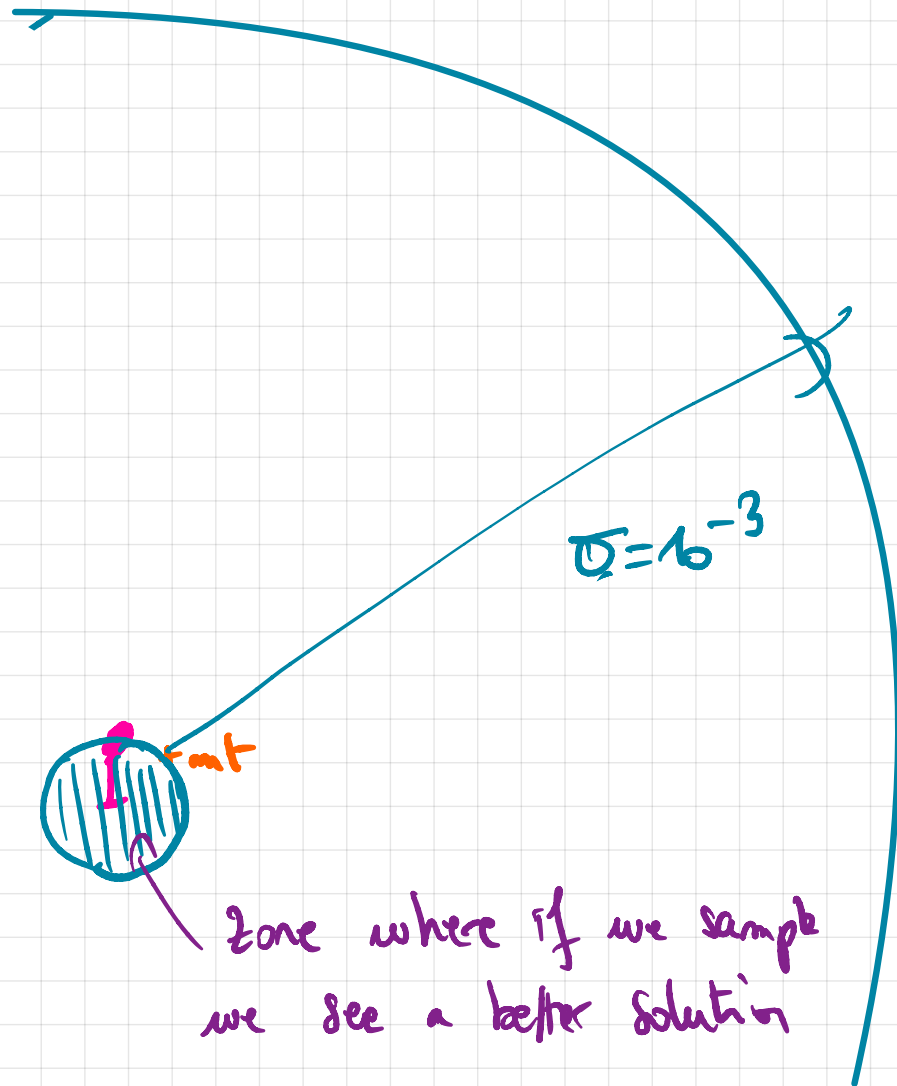
f

Phase III

$$\sigma \gg \|m_t\|$$

Probability $\|x_{t+1}\| \leq \|m_t\|$

very small because proba
to sample better solutions small.



→ Progress too slow

↳ decrease step-size