



# A review of optimization methods for the identification of complex models

Rodolphe Le Riche

► **To cite this version:**

Rodolphe Le Riche. A review of optimization methods for the identification of complex models. Doctoral. France. 2014. <cel-01081272>

**HAL Id: cel-01081272**

**<https://hal.archives-ouvertes.fr/cel-01081272>**

Submitted on 7 Nov 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# **A review of optimization methods for the identification of complex models**

Rodolphe Le Riche  
CNRS and Ecole des Mines de Saint-Etienne

presentation given as part of the MEXICO meeting on  
optimization, La Rochelle, Fr, Nov. 2014  
<http://reseau-mexico.fr>

---

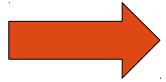
# Goal of the presentation

---

- Why numerical optimization may be useful to the scientist making models ?
  - Provide some concepts on numerical optimization from various backgrounds.
  - Cite a few well-known algorithms, but no claim to be exhaustive.
  - Optimization beyond the algorithm run.
-

# Talk outline

---



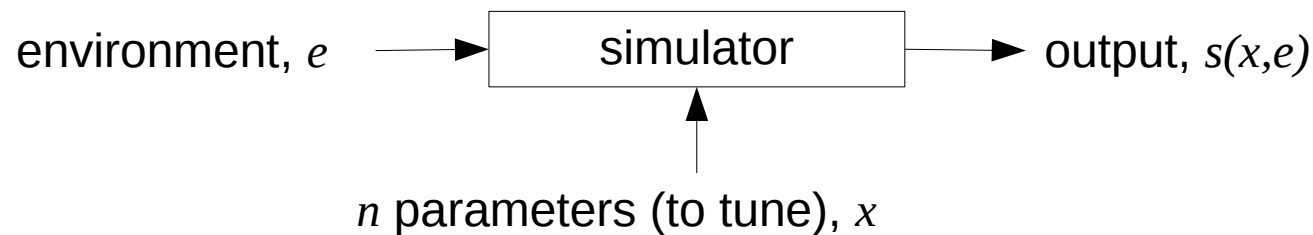
- Why numerical optimization may be useful to scientists making models ?
  - Formulations and problems classification
  - Optimization engines
  - A few references
  - Levenberg-Marquardt algorithm
  - CMA-ES
  - Concluding remarks
  - (other algorithms)
-

# Why numerical optimization may be useful to scientists making models ?

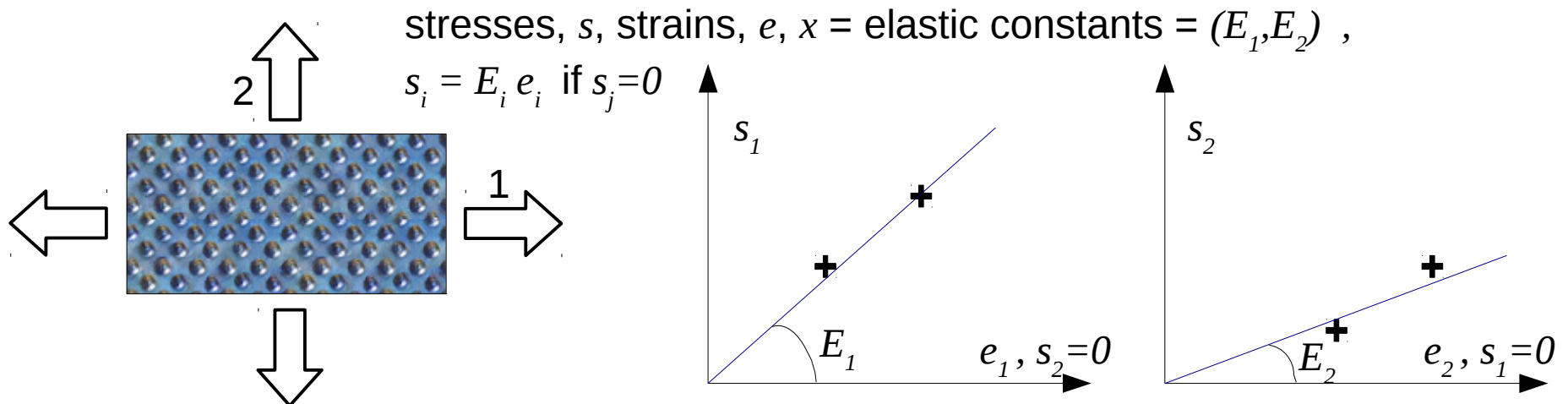
---

## Context of model parameter identification

Ideally, scientists know the physical meaning of the parameters that are in the models and know how to tune them based on specific experiments



1 parameter directly read with 1 experiment, separation of parameter effects  
Example in mechanics, 2 elastic constants to set = slopes of stress-strain curves

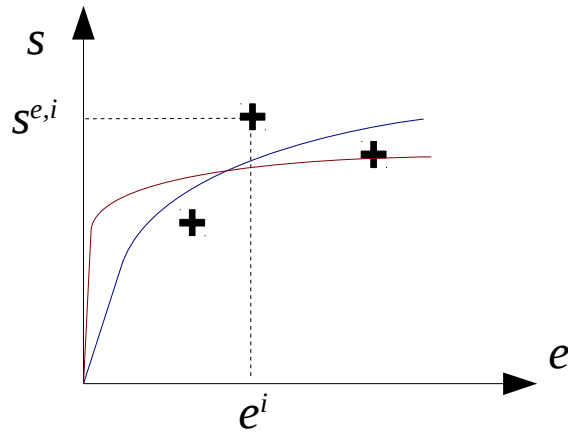


# Why numerical optimization may be useful to models ?

## Context of model parameter identification

---

But the effects of parameters are often difficult to separate



Example : behavior with added contributions  
(linear and non linear stresses)

$$s = Ee + Ke^p, \quad x = (E, K, p)$$

→ formulate an optimization problem : minimize a distance,  $f$ , between experiments and simulations by changing  $x$

$$\min_{x \in S} f(x)$$

$$\text{Classically, } f(x) = \frac{1}{2} \sum_{i=1}^m (s^{e,i} - s(x, e^i))^2 = \frac{1}{2} \|s^e - s(x)\|^2$$

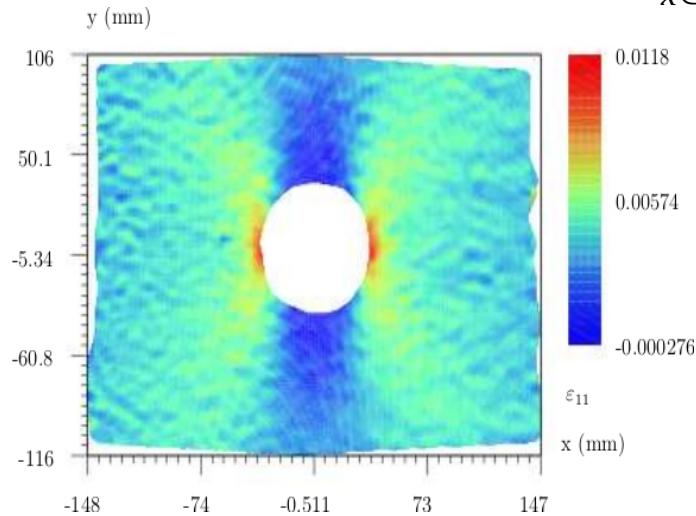
---

# Why numerical optimization may be useful to models ?

## Context of model parameter identification

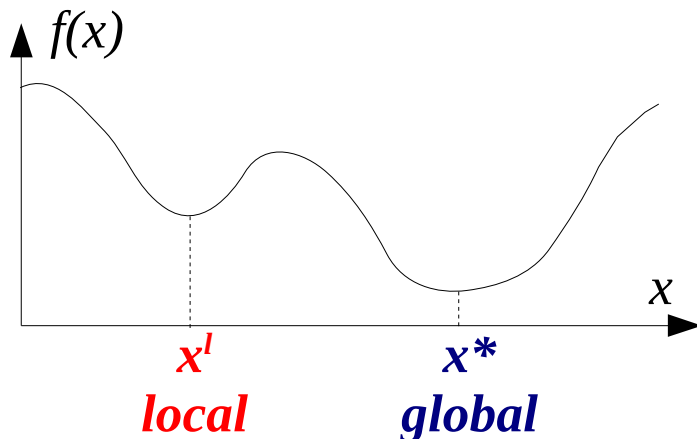
---

To try to solve  $\min_{x \in S} f(x)$  a computer will be helpful



1) because  $S$  will typically be a large space, e.g.,  $R^n$  with  $n \geq 2$

Expl : Identification of 4 elastic constants with 1 experiment (Silva et al., *Identification of material properties using FEMU : application to the open hole tensile test*, Appl. Mech & Matls, 2007)



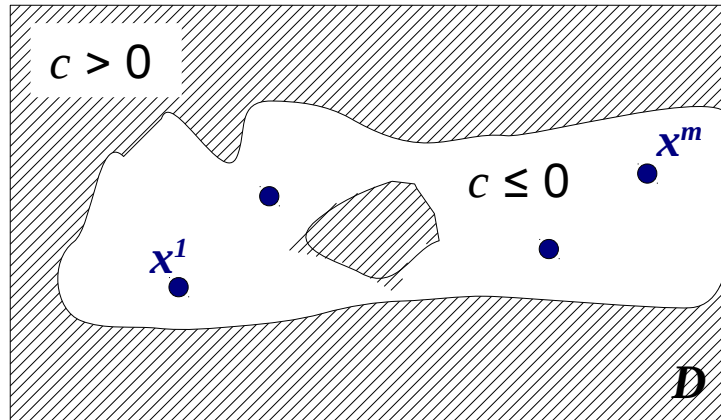
2) because understanding where  $f$  is smallest may be complex, i.e., beyond intuition, because of multimodality.

Expl : Identification of a material non linear law, cf. red and blue solutions on previous slide (+ Le Riche and Guyon, *Least Squares Parameter Estimation and the Levenberg-Marquardt Algorithm : Deterministic Analysis, Sensitivities and Numerical Experiments*, TR, 1999)

---

# Why numerical optimization may be useful to models ?

## Context of DoE or optimal placement



Find the location of points (experiments, sensors, actuators, ... ) that offers best covering while satisfying some constraints (e.g., bounds on parameters, sensors outside protected regions, ... )

$$\min_{x \in S} f(x)$$

$$g(x) \leq 0$$

where

$$S \equiv D^m \subset \mathbb{R}^{m \times d}$$

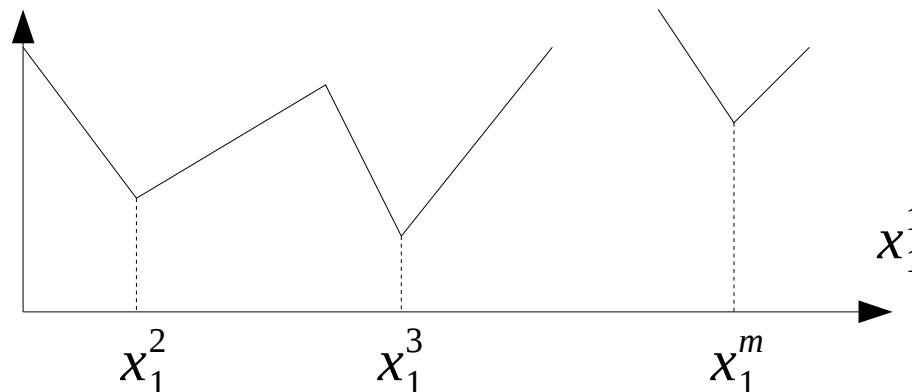
$$x = (x_1^1, \dots, x_d^1, x_1^2, \dots, x_d^m)$$

$$f(x) = -\min_{i \neq j} \text{dist}(x^i, x^j)$$

$$g(x) = \max_{i=1, m} c(x^i)$$

This family of problems is typically high dimensional and multi-modal, not well-suited to human intuition

$$\min_{i \neq 1} \text{dist}(x^1, x^i)$$





# Talk outline

---

- Why numerical optimization may be useful to scientists making models ?



- Formulations and problems classification
  - Optimization engines
  - A few references
  - Levenberg-Marquardt algorithm
  - CMA-ES
  - Concluding remarks
  - (other algorithms)
-

# Formulations

$x$  ,  $n$  optimization variables

$S$  , search space ,  $x \in S$

$f$  , objective or cost function to minimize,  $f : S \rightarrow \mathbb{R}$

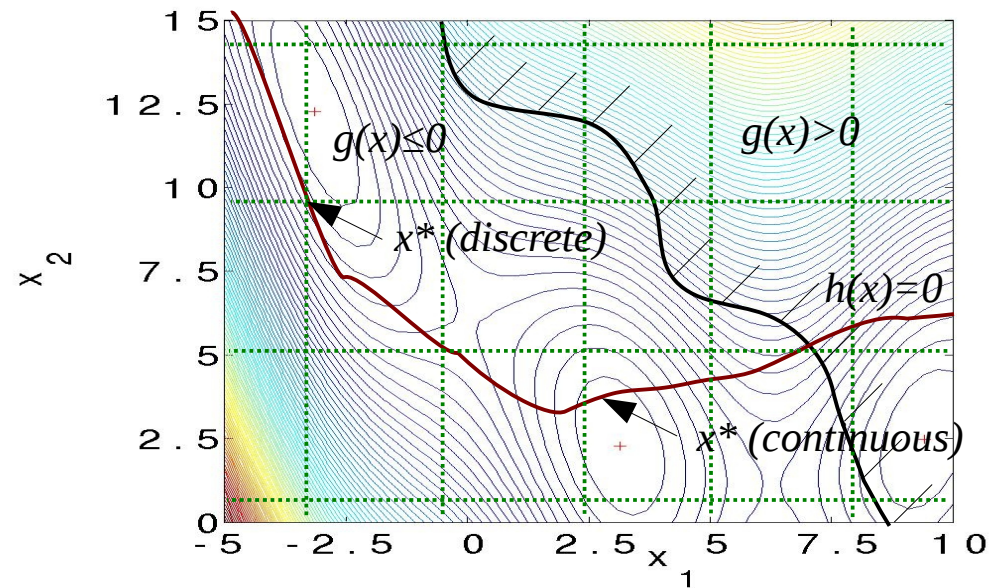
$g$  , inequality constraint,  $g(x) \leq 0$

$h$  , equality constraint,  $h(x) = 0$

$$\min_{x \in S} f(x)$$

$$\text{such that } \begin{aligned} g(x) &\leq 0 \\ h(x) &= 0 \end{aligned}$$

!!! the choice of  $x$  affects the mathematical properties of  $f$ ,  $g$ ,  $h$  and orients the selection of opt. algorithm

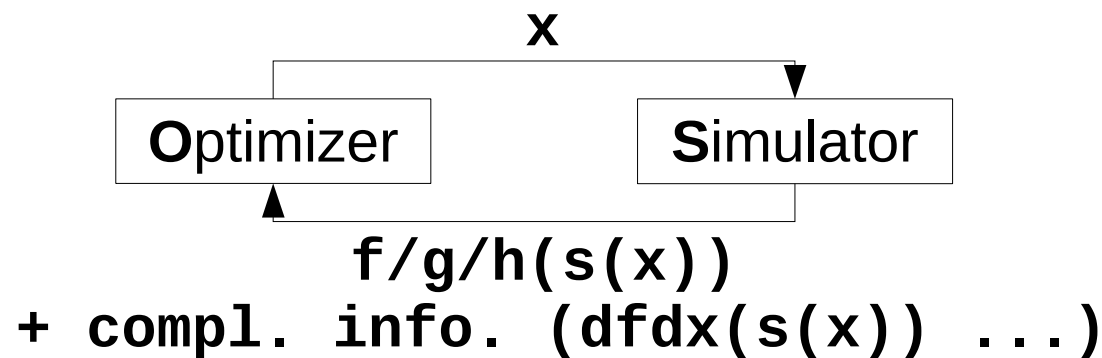
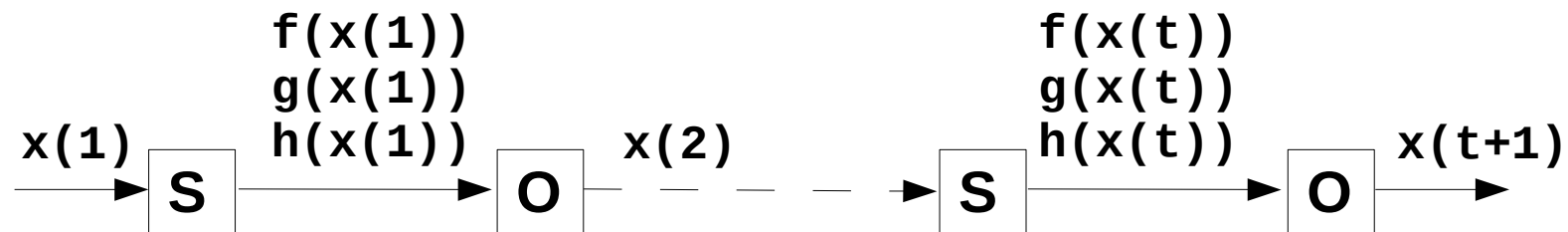


Expl : Displacements ( $f$ ) are non linear functions of thicknesses ( $x = \text{thickness}$ ) but linear functions of  $x = 1/\text{thickness}$  (for statistically determinate structures)

# Optimization algorithms

---

An optimizer is an algorithm that iteratively proposes new  $x$ 's based on past trials in order to approximate the solution to the optimization problem

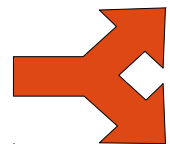



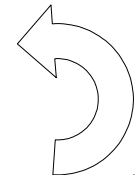
The **cost** of the optimization is the number of calls to the simulator  $s$  (usually = number of calls to  $f$ )

---

# Problems classification

compose the following switches to find the pb. category and the algo

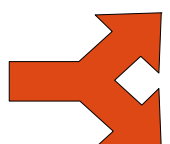
 unconstrained (no  $g$  and  $h$ )  
 constrained



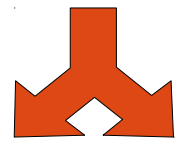
Lagrangians, penalizations, e.g.,  

$$\min_{x \in S} f(x) + p \times \max^2(0, g(x))$$

 mono-objective,  $\min_{x \in S} f(x)$

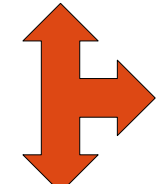
 multi-objective,  $\min_{x \in S} f_1(x) \& \dots f_m(x)$

MO pbs typically have set of (Pareto) solutions, cf. Dimo's talk

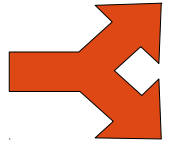


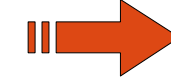
uni-modal of which convex, of which linear → local optimization  
 multi-modal → global optimization

continuous variables  
 $S \subset \mathbb{R}^n$

 integer programming  
 $S \subset \mathbb{N}^n$  no gradients  $x(\dots)$


mixed programming  
 $S \subset \{\mathbb{R}^{n1} \cup \mathbb{N}^{n2}\}$

 costly simulations → metamodels  
 non costly simulations costly = sim. budget < 1000. cf. victor's talk

 number of variables  $n$ , cost of optimizer, ...

# Talk outline

---

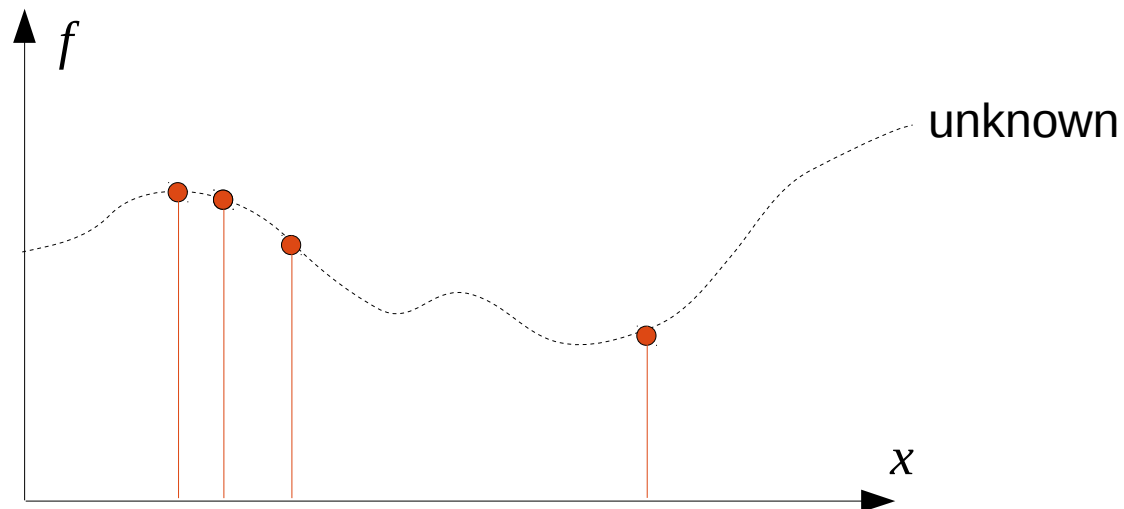
- Why numerical optimization may be useful to scientists making models ?
  - Formulations and problems classification
  -  • Optimization engines
  - A few references
  - Levenberg-Marquardt algorithm
  - CMA-ES
  - Concluding remarks
  - (other algorithms)
-

# Optimization engines

---

$$\min_{x \in S} f(x)$$

and the function's landscape is known through (costly) point measures



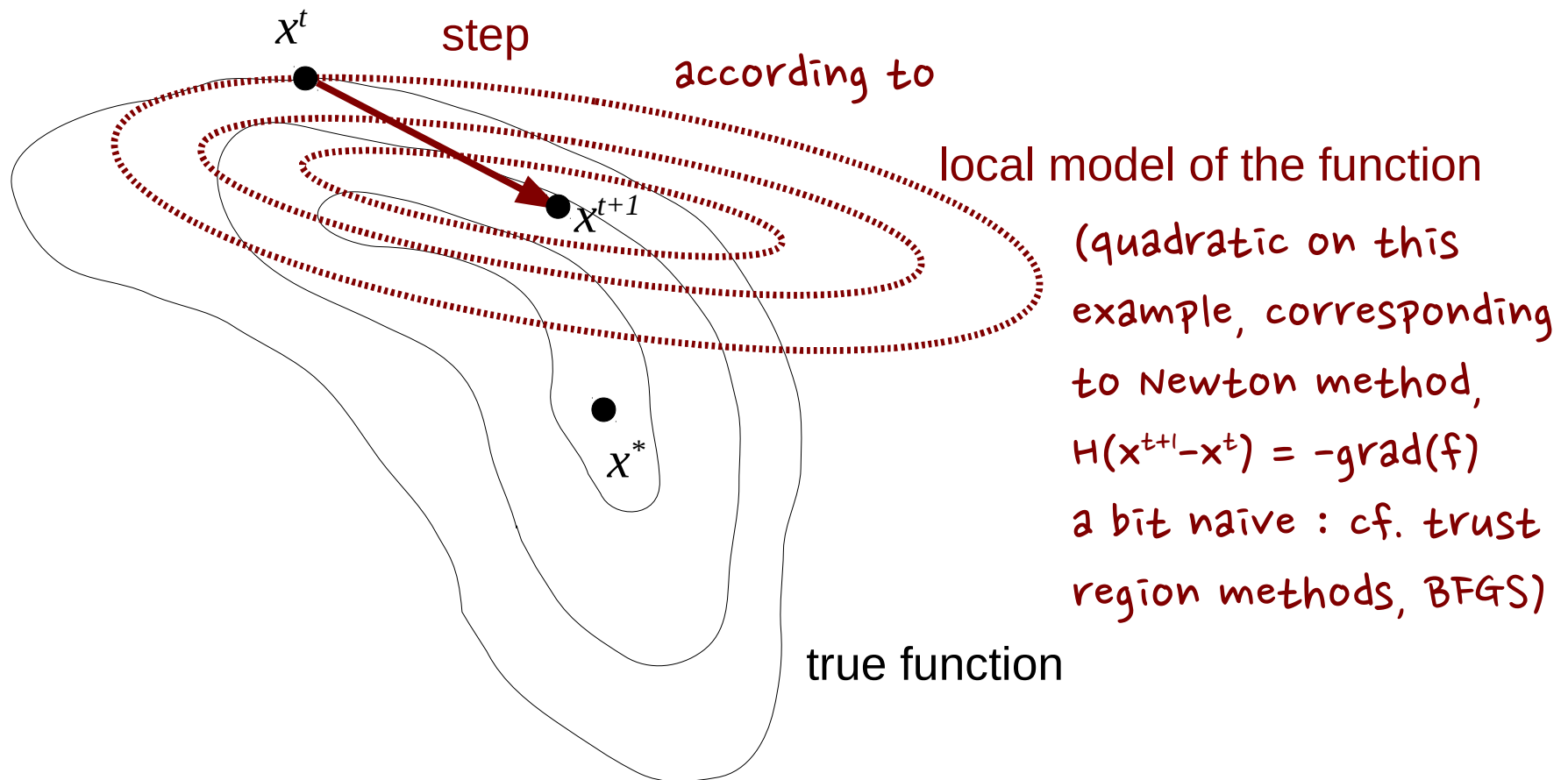
Optimization engine = basic principle to minimise the function

---

# Optimization engines

## Local model of the function

---

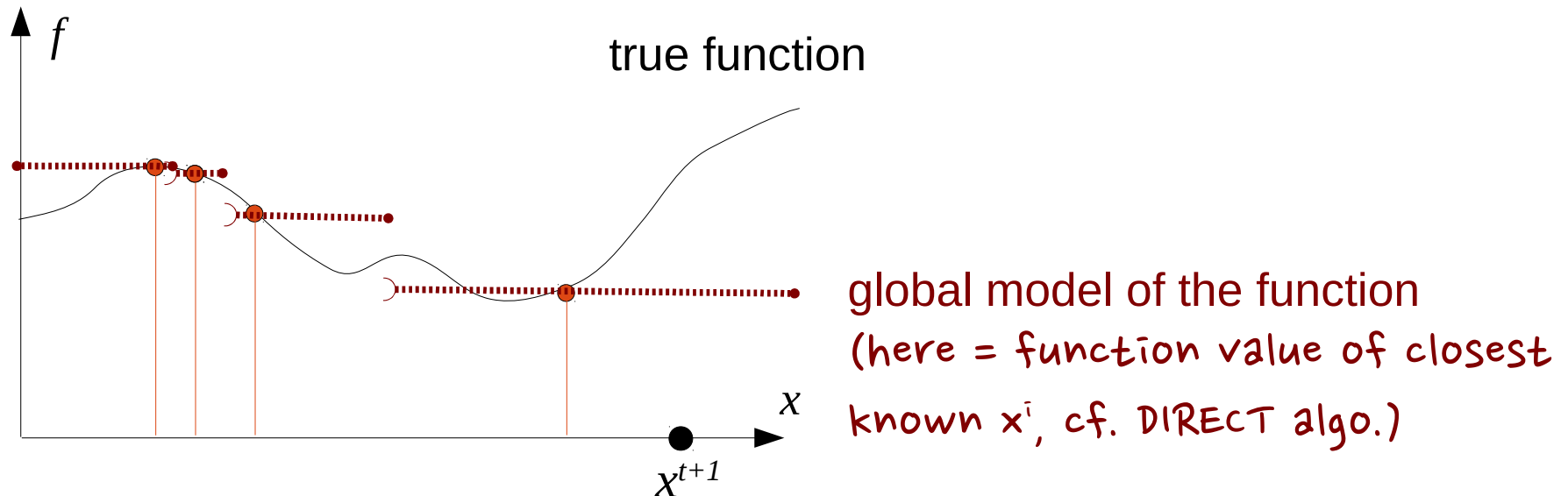


- For local optimization.
  - Efficient (low number of calls to  $f$ ) because focuses on neighborhood of current iterate, suitable for  $n$  large ( $>1000$ ).
  - At the risk of missing global optimum if  $f$  multi-modal.
-

# Optimization engines

## Global model of the function

---



The next iterates  $x^{t+1}$  are compromises between intensification (search in low areas) and exploration (search in unknown areas).

For global optimization.

Efficiency depends on the agreement between the model and the true function but searching the whole  $S$  is subject to the « curse of dimensionality ».

---

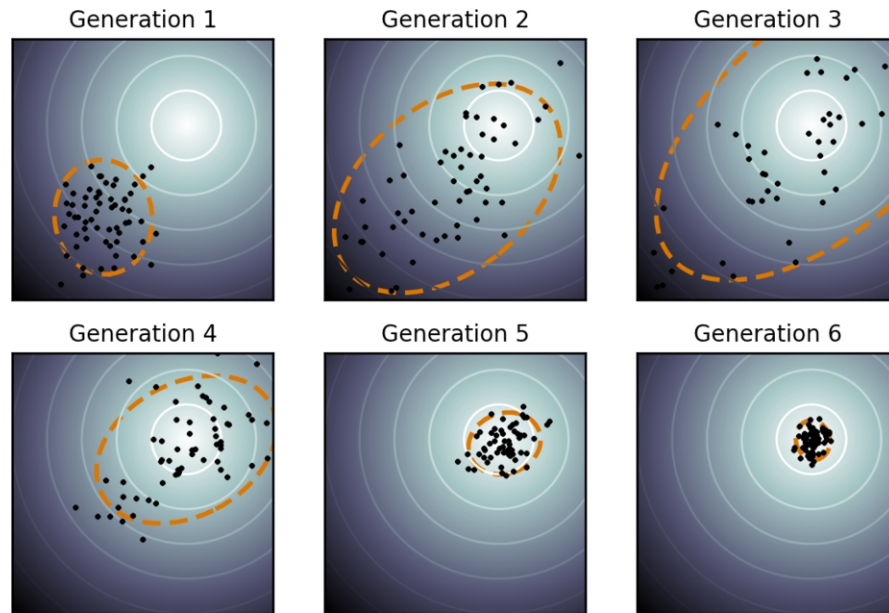


# Optimization engines

## local sampling

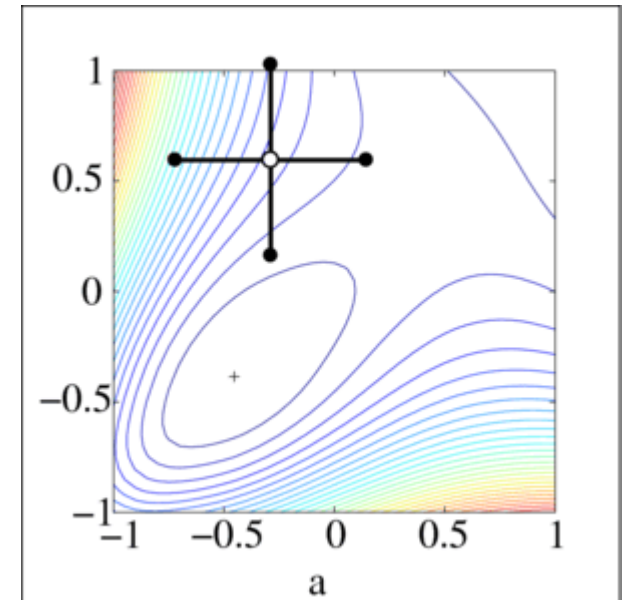
### sampling in probability

e.g., from a multi-normal density, cf. CMA-ES algorithm



### sampling with a geometrical pattern

e.g., from a pattern search algorithm



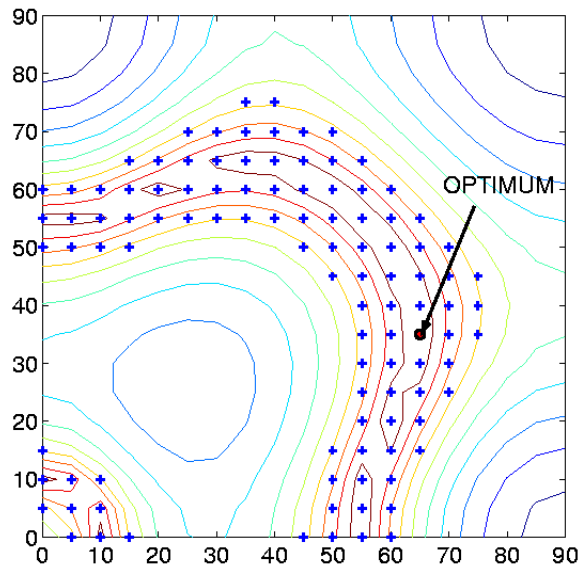
- Sampling from a current point with the notion of step.
- Sampling → methods do not require  $f$  continuity or regularity.
- Sampling in probability →  $S$  can be continuous, or discrete or mixed.
- May converge locally.
- Efficient and more suitable to large  $n$  ( $>10$ ) than global sampling methods.

(both graphs from wikipedia, cma-es – left – and pattern search – right –)

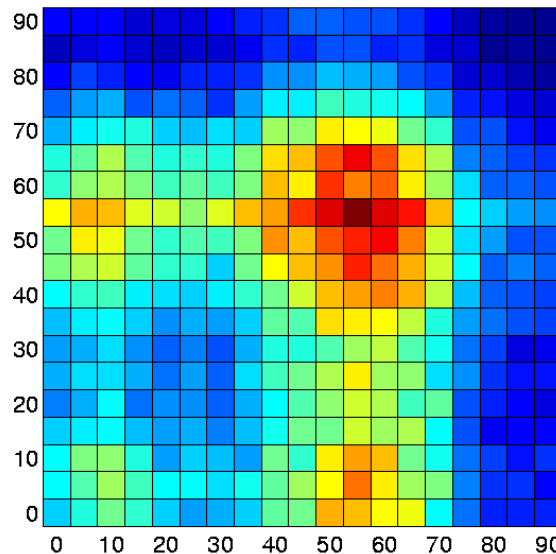
# Optimization engines global sampling

The next  $x$  iterates can be sampled anywhere in  $S$  from a density learned on observed points. Expl. Estimation of Distribution Algorithms (illustration below), Branch and Bound, ... .

$f(x)$  and selected points



$$p_X(x) = \prod_{i=1}^n p_i(x_i)$$




(from Grosset et al., *A double distribution statistical algorithm for composite laminate optimization*, SMO, 2006)

- Sampling  $\rightarrow$  methods do not require  $f$  continuity or regularity.
- Sampling in probability  $\rightarrow S$  can be continuous, or discrete or mixed.
- May converge globally.
- Typically costly.

# Talk outline

---

- Why numerical optimization may be useful to scientists making models ?
  - Formulations and problems classification
  - Optimization engines
  -  • A few references
  - Levenberg-Marquardt algorithm
  - CMA-ES
  - Concluding remarks
  - (other algorithms)
-

# Optimization engines the quest for the best optimizer

---

4 « optimization engines » and hundreds of algorithms that vary in their implementation and composition of these principles.

What is the best optimizer ?

If you have no a priori on the functions you want to optimize, there is no best optimizer (*No Free Lunch Theorem in optimization*, D. Wolpert, 1997).

But there are always a priori's on the considered  $f$ 's (no « needle in the haystack », class of problems such as identification of a given physical model, ... ).

→ The best algorithm depends on the considered  $f$ 's and the allocated computational budget.

---

# A few references (1)

---

No claim to be exhaustive, fair, well-informed, ...  
just a few references to be practical.

$S \subset \mathbb{R}^n$  , **continuous optimization**

**Constrained linear programming**

→ *Simplex algorithm*, Dantzig, 1947.

$$\min_{x \in \mathbb{R}^n} c^T x \text{ such that } Ax - b \leq 0$$

Try to achieve such formulation because  $n$  and the number of constraints can be very large (thousands) and the resolution almost immediate.

**Unconstrained, uni-modal, smooth function**

→ *BFGS algorithm*, (Broyden-Fletcher-Goldfarb-Shanno, 1970)

A quasi-Newton method requiring the calculation of gradient of  $f$ .

→ *Levenberg-Marquardt* (Levenberg 1944, Marquardt 63)

A quasi-Newton method specialized for non-linear least-squares (model identification!). Gradient of  $f$  to calculate.

**Unconstrained, moderately multi-modal function without derivative**

→ *NEWUOA algorithm* (Powell, 2004)

A trust-region algo. with quadratic approximation, no function derivative, sparse sampling ( $\sim 2n+1$  interpolation points),  $n$  up to 160 variables.

---

# A few references (2)

---

## $S \subset \mathbb{R}^n$ , continuous optimization

(– Continued – **unconstrained, moderately multi-modal function without derivative**)

→ *Nelder-Mead* algorithm (1965), robust and simple, may stop on non-stationary points, up to  $n \sim 30$ . Cf. *A constrained, globalized and bounded Nelder-Mead algorithm for eng. optim.*, Luersen et al., SMO 2003.

**Unconstrained, multi-modal function without derivative**

→ *CMA-ES* algorithm (Hansen et al., from 1996 on)  
may be the *BIPOP-CMA-ES* variant for restarts. up to  $n \sim 100-1000$ .

**Unconstrained, multi-modal function without derivative and expensive functions**

→ *CMA-ES* with surrogate for moderately expensive functions (may be the *saACM-ES* variant by Loshchilov et al., 2013)  
→ *EGO* (Jones, 1998) for very expensive functions (budget  $< 1000$ ).

**Constrained non linear continuous optimization**

→ *Sequential Quadratic Programming*  
a trust region iterative method with a series of Quadratic Problems solved.  
Up to  $n \sim 100-1000$ .

---

# A few references (3)

---

$S \subset \mathbb{N}^n$  , integer programming

**linear integer programming**

$$\min_{x \in \mathbb{Z}^n} c^T x \text{ such that } Ax - b \leq 0$$

→ *cutting plane* or *branch and bound* methods

solutions in polynomial time, very large number of variables and constraints possible. As in continuous linear programming, work on the formulation to obtain such a problem, if possible.

Cf. *Introduction to operations research*, Hillier & Lieberman, 1990.

**non linear integer programming**

→ Univariate Marginal Density Algorithm, *UMDA*, Baluja 1994 as PBIL, Mühlenbein 1996. For its simplicity.

→ *MIMIC* ( Mutual Information Maximizing Input Clustering ) algorithm : De Bonnet, Isbell and Viola, 1997. Accounts for variables couplings.

→ other problem specific metaheuristics in particular for combinatorial problems, e.g., *Evolutionary Algorithms*, (cf. *Introduction to evolutionary computing*, Eiben and Smith, 2003).

EAs = Possible line of approach for mixed problems,  $S \subset \{\mathbb{R}^{n1} \cup \mathbb{N}^{n2}\}$

All these algorithms are expensive (1000-10<sup>6</sup> calls to  $f$ ).

---

# A few references (4)

---

**Constrained mixed non linear optimization** ,  $S \subset \mathbb{R}^{n1} \cup \mathbb{N}^{n2}$

**box constraints**  $x^{\min} \leq x \leq x^{\max}$

are simple to handle inside the algorithms, either by satisfying the Karush, Kuhn and Tucker optimality conditions, or by re-sampling  $\rightarrow$  look for versions of the above algorithms with box constraints, e.g., *BFGS-B*.

**non linear equality constraints**

can be turned into two inequality constraints

$$h(x)=0 \Leftrightarrow h(x) \leq 0 \text{ and } -h(x) \leq 0$$

**non linear inequality constraints**

can be merged into the objective function by a penalty technique, a recommended form being the *augmented Lagrangian* (Rockaffelar, 1973)

$$L_{\text{aug}}(x, \lambda) = \begin{cases} f(x) - \lambda^2 / 4r & \text{if } g(x) \leq -\frac{\lambda}{2r} \\ f(x) + \lambda g(x) + r [g(x)]^2 & \text{else} \end{cases}$$

(index  $\lambda$  and  $g$  with  $i$  and sum if many constraints )


$\lambda =$  Lag. multipliers,  $\min_{x \in S} L_{\text{aug}}(x, \lambda^t) \rightarrow x(\lambda^t)$  ,  $\lambda^{t+1} \leftarrow \lambda^t + 2r g(x(\lambda^t))$

---



# Talk outline

---

- Why numerical optimization may be useful to scientists making models ?
  - Formulations and problems classification
  - Optimization engines
  - A few references
  -  • Levenberg-Marquardt algorithm
  - CMA-ES
  - Concluding remarks
  - (other algorithms)
-

# Non-Linear Least Squares and the Levenberg-Marquardt algorithm

---

$$\min_{x \in S \subset \mathbb{R}^n} f(x) \quad \text{where} \quad f(x) = \frac{1}{2} \sum_{i=1}^m (s^{e,i} - s(x, e^i))^2 = \frac{1}{2} \|s^e - s(x)\|^2$$

Motivation : if  $s^{e,i} = s_i(x^*) + \epsilon_i$  where  $\epsilon_i \sim N(0, \sigma_i^2)$  ,  
 $\max$  likelihood measures  $\Leftrightarrow \min_{x \in S} f(x)$

Principle : iteratively solve,

$$\left\{ \begin{array}{l} \min_{x \in S \subset \mathbb{R}^n} \frac{1}{2} \|s(x^k) + \nabla s(x^k)(x - x^k) - s^e\|^2 \\ \text{such that } \frac{1}{2} \|x - x^k\|^2 \leq \frac{1}{2} \delta_k^2 \end{array} \right.$$

a local quadratic approximation with trust region

Refs. :

Le Riche and Guyon, *Least Squares Parameter Estimation and the Levenberg-Marquardt Algorithm : Deterministic Analysis, Sensitivities and Numerical Experiments*, TR, 1999.

Besson et al., *Object Oriented Programming applied to the Finite Element Method - Part II : Application to Material Behavior*, Europ. J. of Finite Elements, 1998.

---

# Non-Linear Least Squares and the Levenberg-Marquardt algorithm

---

Optimality conditions (gradient of Lagrangian cancels out) yield :

0. Initialize :  $\nu > 1$  ,  $x^0$  ,  $\lambda_0 > 0$  ,  $\Delta^{\min} > 0$
1. 
$$\left( \nabla s(x^k)^T \nabla s(x^k) + \lambda_k I \right) (x^{k+1} - x^k) = - \nabla s(x^k)^T (s(x^k) - s^e)$$
$$\Rightarrow x^{k+1}$$
2. If  $f(x^{k+1}) < f(x^k)$  ,  $\lambda_k \leftarrow \lambda_k / \nu$  ,  $k \leftarrow k+1$  , goto 1  
else if  $\|x^{k+1} - x^k\| > \Delta^{\min}$  ,  $\lambda_k \leftarrow \lambda_k \times \nu$  , goto 1  
else stop

- LM is specialized for least squares.
  - LM is based on a local model of the function. cf. "optim. engines"
  - Requires the gradient of the simulator,  $s$ , w.r.t.  $x$ .
  - $n$  up to [100-1000] , but if the gradient is estimated by finite differences ( $n+1$  calls to  $s$  each time) it becomes costly.
-

# Non-Linear Least Squares and the Levenberg-Marquardt algorithm

---

- $\nabla f(x) = \nabla s(x)^T (s(x) - s^e)$

when progress is rare,  $\lambda_k \rightarrow$  large and LM becomes a steepest descent search, with guaranteed progress at the price of small steps.

- $\nabla^2 f(x) = \nabla s(x)^T \nabla s(x) + \sum_{i=1}^m \nabla^2 s_i(x) (s_i(x) - s_i^e) \xrightarrow{s \rightarrow s^e} \nabla s(x)^T \nabla s(x)$

i.e., LM algorithm is a quasi-Newton method near the solution when  $\lambda_k \rightarrow 0$ .

- In case of non identifiability, LM performs a smallest dist. to  $x^0$  regularization (rigorous when  $s$  linear in  $x$ ),

$$\left\{ \begin{array}{l} \min_{x \in S} \frac{1}{2} \|x - x^0\|^2 \\ \text{such that } f(x) = \min_{y \in S} f(y) \end{array} \right.$$

---

# Levenberg-Marquardt post-processing local sensitivity analysis

---

$$s(x + \Delta x) = s(x) + \nabla s(x) \Delta x + O(\Delta x^2)$$

$$\|s(x + \Delta x) - s(x)\|^2 = \Delta x^T \underbrace{\nabla s(x)^T \nabla s(x)}_{n \times n} \Delta x + O(\Delta x^3)$$

$(\kappa_i^2, v_i)$  singular values / vectors

$$\frac{\|s(x + \alpha v_i) - s(x)\|^2}{\|\alpha v_i\|^2} = \kappa_i^2 + O(\alpha) \quad , \quad 0 < \alpha \ll 1$$

$\Rightarrow \kappa_i^2$  sensitivities to parameters

Sensitivity to experiments  $s^e$

$$\approx \text{cond}_2(\nabla s(x)^T \nabla s(x)) = \frac{\max_i \kappa_i^2}{\min_i \kappa_i^2}$$

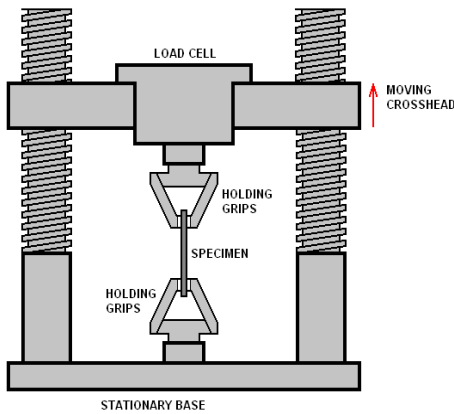
---

# Levenberg-Marquardt algorithm

## Expl. Identification of material model

Identify 4 non-linear material parameters

2 experiments



$$x = (E, K, n, R_i)$$

$$e = \{\epsilon_1, \dots, \epsilon_m\}$$

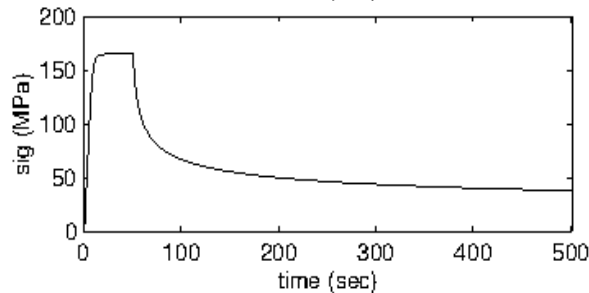
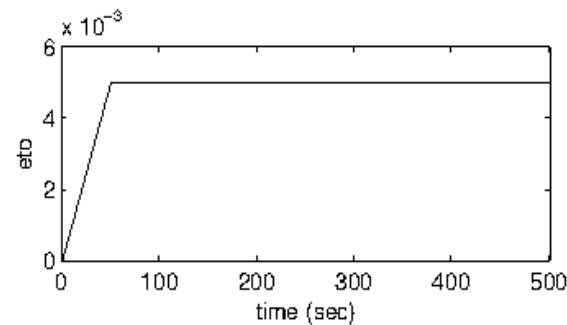
$$s(x) = \{s(x, \epsilon_1), \dots, s(x, \epsilon_m)\}$$

simulator

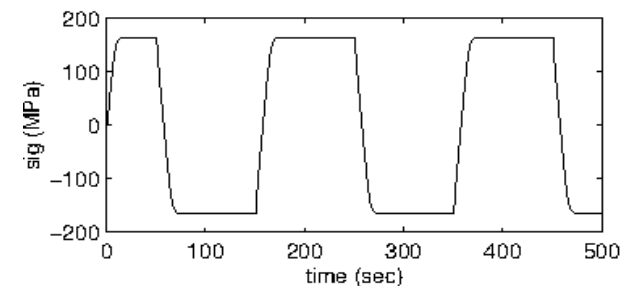
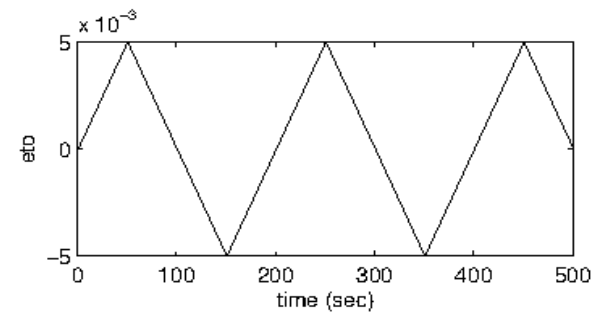
$$\begin{cases} \dot{p} = \eta \left( \frac{|\sigma| - R_i}{K} \right)^n \\ \dot{\epsilon}_i = \eta \dot{p} \\ \dot{\epsilon}_e = \dot{\epsilon} - \dot{\epsilon}_i \end{cases}$$

and  $\begin{cases} \sigma = E \epsilon_e \\ \eta = \text{sign}(\dot{\epsilon}) \end{cases}$

relaxation



fatigue



# Levenberg-Marquardt algorithm

## Expl. Identification of material model

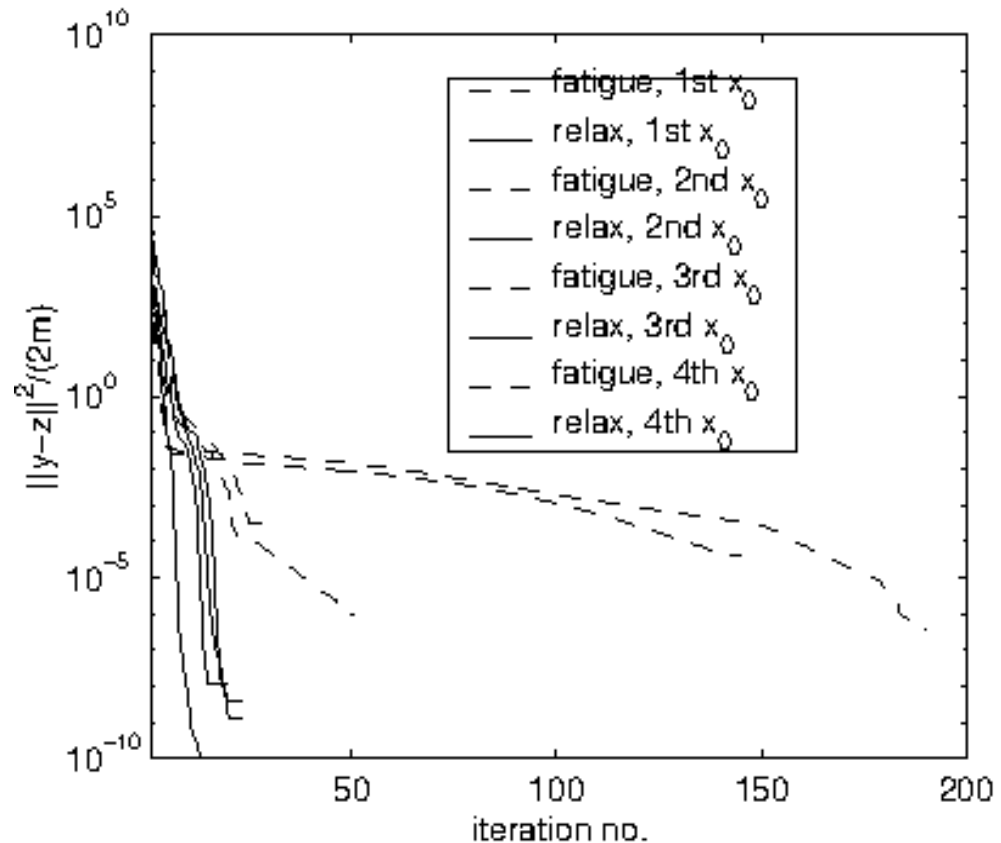
---

It is easier to identify the model with the relaxation experiment.

This is seen through the condition number ( $\sim 10^4$  for the relaxation against  $\sim 10^8$  for the fatigue)

The eigenvector associated to the smallest eigenvalue of the fatigue experiment,  $v_4$ , corresponds to the asymptote

$$s(x^*) = \text{const}$$
$$v_4^T \cdot \nabla s(x^*) \approx 0$$



# From Levenberg-Marquardt to evolution strategies

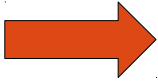
---

- $s$  non-linear ,  $n > m \Rightarrow$  local optima
  - other, non smooth, distances ,  $f_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$   
for example  $f_1$  for sparsity
  - constraints handling through penalties  
 $\Rightarrow$  a more general optimization algorithm discussed next
-



# Talk outline

---

- Why numerical optimization may be useful to scientists making models ?
  - Formulations and problems classification
  - Optimization engines
  - A few references
  - Levenberg-Marquardt algorithm
  -  • CMA-ES
  - Concluding remarks
  - (other algorithms)
-

# Stochastic optimizer : CMA-ES

---

(N. Hansen et al., from 1996 on, then co-work with A. Auger, M. Schoenauer, M. Sebag, D. Brockhoff, I. Loshchilov)

CMA-ES = *Covariance Matrix Adaptation Evolution Strategy* = optimization through multi-normal law sampling and updating.

Optimization engine = local sampling in probability.

A state of the art method for stochastic optimization.

---

# Simplified flow chart of *ES-(1+1)*

---

Initializations :  $x, f(x), m, C, t_{max}$ .

While  $t < t_{max}$  do,

    Sample  $N(m, C) \rightarrow x'$

    Calculate  $f(x'), t = t + 1$

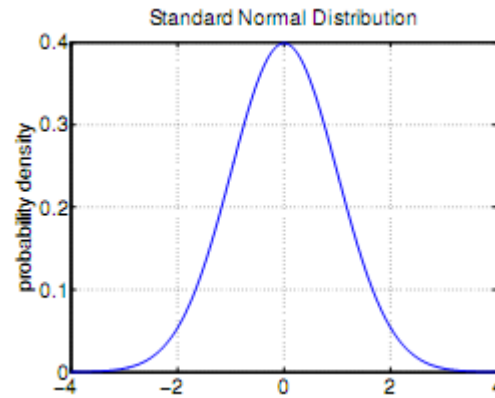
    If  $f(x') < f(x), x = x', f(x) = f(x')$  Endif

    Update  $m$  (e.g.,  $m = x$ ) and  $C$

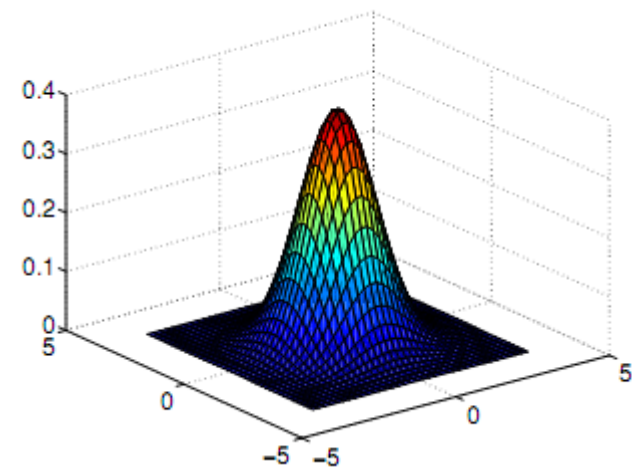
End While

Normal law

$N(m, C)$



2-D Normal Distribution



# Numerical sampling of a multi-normal law

To sample  $Y \sim N(\mathbf{m}, \mathbf{C})$ , perform the eigen analysis

$$\mathbf{C} = \mathbf{U} \mathbf{D}^2 \mathbf{U}^T \quad \text{and then notice} \quad \mathbf{Y} = \mathbf{m} + \mathbf{U} \mathbf{D} \mathbf{E}, \quad \mathbf{E} \sim N(0, \mathbf{I})$$

( prove that  $EY = \boldsymbol{\mu}$  and  $\text{Cov}(Y) = \mathbf{C}$ ,  
+ a Gaussian vector is fully determined by its 2 first moments )

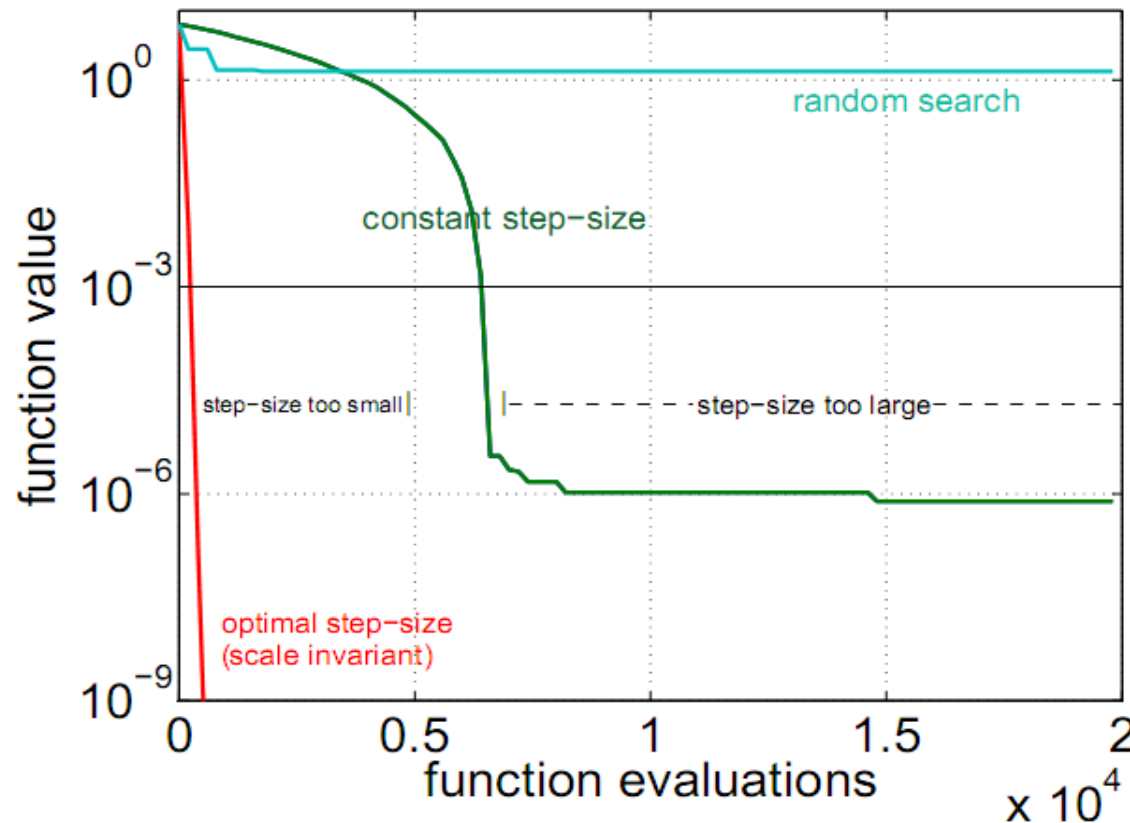
**In R,**

```
Ceig <- eigen(C)
y <- mu + Ceig$vector %*% diag(sqrt(Ceig$value))
      %*% matrix(rnorm(n))
```



# Adapting the step size ( $C$ ) is important

---



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

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

(A. Auger et N.  
Hansen, 2008)

Above, ES(1+1) with isotropic step :  $\mathbf{C} = \sigma^2 \mathbf{I}$  ,  $\sigma$  is the step size.  
With an optimal step size (  $\approx \|\mathbf{x}\|/n$  ) for the sphere function, performance  
only degrades in  $O(n)$  !

---

# Simplified CMA-ES flow chart

CMA-ES is an evolution strategy  $ES-(\mu, \lambda)$  :

Initializations :  $m, C, t_{max}, \mu, \lambda$

While  $t < t_{max}$  do,

    Sample  $N(m, C) \rightarrow x^1, \dots, x^\lambda$

    Calculate  $f(x^1), \dots, f(x^\lambda)$  ,  $t = t + \lambda$

    Rank :  $f(x^{1:\lambda}), \dots, f(x^{\lambda:\lambda})$

    Update  $m$  and  $C$  with the  $\mu$  best,  $x^{1:\lambda}, \dots, x^{\mu:\lambda}$

End While

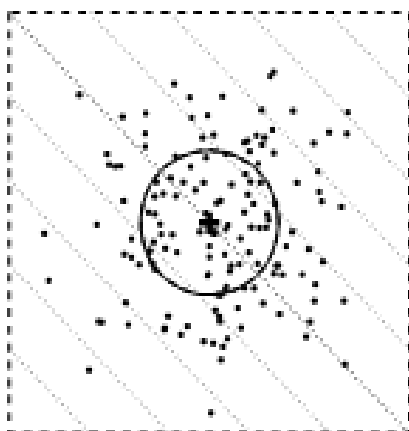
$m$  and  $C$  are updated with

- the best **steps**,
  - a **time cumulation** of these steps.
-

# Simplified CMA-ES : $C$ adaptation with the last good steps

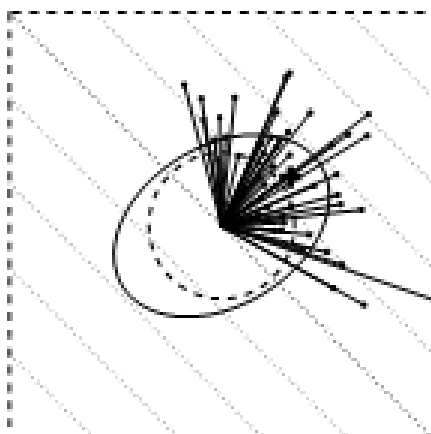
(A. Auger et N. Hansen, 2008)

Initialization :  $m \in S$  ,  $C = I$  ,  $c_{cov} \approx 2/n^2$



sampling

$$\begin{aligned}x^i &= m + y^i \\y^i &\propto N(0, C) \\i &= 1, \dots, \lambda\end{aligned}$$

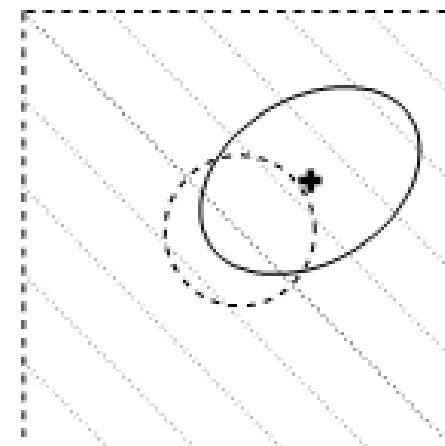


selection

$$y_w = \frac{1}{\mu} \sum_{i=1}^{\mu} y^{i:\lambda}$$

$C$  rank 1 update

$$C \leftarrow (1 - c_{cov})C + c_{cov} \mu y_w y_w^T$$

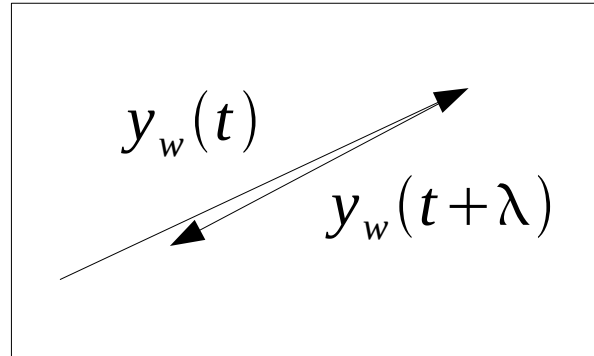


update  $m$

$$m \leftarrow m + y_w$$

# Simplified CMA-ES : time cumulation of good steps

---



When successive steps are anti-correlated, the step size should decrease, which is impossible so far since  $y_w y_w^T = (-y_w) (-y_w)^T$ .

Cumulation (time exponential damping) :

$$p_c \leftarrow \underbrace{(1-c_c)}_{\text{damping}} p_c + \underbrace{\sqrt{1-(1-c_c)^2} \sqrt{\mu}}_{\text{normalization}} y_w$$

$$C \leftarrow (1-c_{cov}) C + c_{cov} p_c p_c^T \longleftarrow \text{rank 1}$$

$$c_c \approx 3/n \quad , \quad c_{cov} \approx 2/n^2$$



# Further comments on CMA-ES

---

(check the BBOB web site for latest versions)

## Invariance properties

- CMA-ES invariant w.r.t. monotonic transformation of  $f$ , e.g.,  $f \leftarrow \log(f)$ , because comparison based.
- CMA-ES invariant w.r.t. change in coordinate system (rotation + translation) because of update formula.

## Additional features of state of the art CMA-ES ( / previous slides) :

- steps are weighted,  $y_w = \sum_{i=1}^{\mu} w_i y^{i:\lambda}$
  - Rank 1 and  $\mu$  updates of  $C$  carried out simultaneously,
  - Global scaling of the covariance matrix,  $C \rightarrow \sigma^2 C$  .
  - Restarts with a small and a large population size (BIPOP version).
-

# Concluding remarks

## Optimization added value beyond the algorithm run

---

- Optimization starts with a precise description of what to do with a simulator (mathematical formulation) in a decision making process.
  - Problem solving will often show simulator / formulation flaws : repeat the optimization.
  - Use the added value of the optimization effort : metamodels are constructed, (local or global) sensitivities become available.
-