



CentraleSupélec

MsC in Data Sciences & Business Analytics

Optimization Course

Gif-sur-Yvette

December 13, 2017

Charles SOUSSEN

`charles.soussen@centralesupelec.fr`

Summary

1. Introduction.
2. Unconstrained optimization.
 - Derivative-free algorithms.
 - Gradient-based algorithms.
 - Least squares algorithms.
3. Least-square problems.
4. Constrained optimization.
5. Global / discrete optimization.

Organization

8 courses:

- Course 1: lecture.
- Courses 2 to 8: 1h30 lecture and 1h30 practical session (matlab).
- Dec. 14, 2017: exam.

1. Introduction

General formulation (unconstrained)

Minimize $f(x)$.

Terminology:

- $x \in \mathbb{R}^n$: **variables**.
- $f : \mathbb{R}^n \mapsto \mathbb{R}$: **objective function** (criterion).

Resolution:

- Existence and unicity of the solution?
- Which numerical algorithm?
- Algorithm complexity: computation time, memory storage.

General formulation (constrained)

$$\text{Minimize } f(x) \quad \text{subject to} \quad \begin{cases} g_i(x) \leq 0, & i = 1, \dots, p \\ h_i(x) = 0, & i = 1, \dots, q. \end{cases}$$

Terminology:

- $x \in \mathbb{R}^n$: variables.
- $f : \mathbb{R}^n \mapsto \mathbb{R}$: objective function.
- $g_i : \mathbb{R} \mapsto \mathbb{R}$: inequality constraints.
- $h_i : \mathbb{R} \mapsto \mathbb{R}$: equality constraints.

Reformulation

Minimize $f(x)$ **subject to** $x \in \mathcal{D}$.

Feasible set: $\mathcal{D} = \{x \in \mathbb{R}^n, g(x) \leq \mathbf{0}_p, h(x) = \mathbf{0}_q\}$

- $g(x) = \begin{bmatrix} g_1(x) \\ \vdots \\ g_p(x) \end{bmatrix}$ gathers the p inequality constraints.
- $h(x) = \begin{bmatrix} h_1(x) \\ \vdots \\ h_q(x) \end{bmatrix}$ gathers the q equality constraints.

Example 1: Knapsack problem

Jo goes hitch hiking. The maximum weight allowed in his knapsack is W . Each article $i = 1, \dots, n$ he can take weight w_i and has a usefulness u_i . What articles should be taken in the knapsack to maximize the usefulness?

Example. $W = 25$, and:

| i | w_i | u_i |
|-----|-------|-------|
| 1 | 25.00 | 40 |
| 2 | 12.50 | 35 |
| 3 | 11.25 | 18 |
| 4 | 5.00 | 4 |
| 5 | 2.50 | 10 |
| 6 | 1.25 | 2 |

Mathematical formulation

1. What is the search space?

$$x \in \mathcal{D} = \{0, 1\}^n$$

$x_i = 1 \Leftrightarrow$ bring article # i in the knapsack

2. What is the objective function?

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

3. What constraints?

$$\sum_{i=1}^n w_i x_i \leq W$$

\Rightarrow *Combinatorial* (discrete) optimization problem.

Example 2: Traveling salesman problem

Given a list of cities and the distances between each pair of cities, what is the shortest possible route that visits each city exactly once and returns to the origin city?

- Search space? \mathcal{D} : set of all permutations of $\{1, \dots, n\}$.

\Rightarrow *Combinatorial* (discrete) optimization problem.

Example 3: Ressource allocation: example

A factory produces two products P1 and P2 and has resources R1 (equipment), R2 (human), R3 (raw materials) in limited quantities.

P1 is sold for 6 euros per unit, P2 is sold for 4 euros per unit.

Example.

| | P1 | P2 | Availability |
|--------|----|----|--------------|
| R1 | 3 | 9 | 81 |
| R2 | 4 | 5 | 55 |
| R3 | 2 | 1 | 20 |
| Profit | 6 | 4 | |

Questions. Mathematical formulation? Feasible domain?

Example 3: Ressource allocation

- Bandwith allocation in wireless communication.
- Strategic planning: allocating scarce resources (e.g., human resources) among the various projects or business units.
- *etc.*

Linear programming problems.

Example 4: Portfolio optimization

Choosing the proportions of various assets to be held in a portfolio, in such a way as to make the portfolio better than any other according to some criterion.

The criterion will combine considerations of:

- the expected value of the portfolio's rate of return
- the return dispersion
- other measures of financial risk.

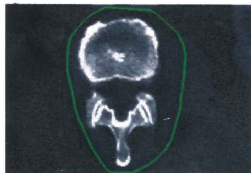
Example 5: Mechanical design of structures

Bridge topology optimization

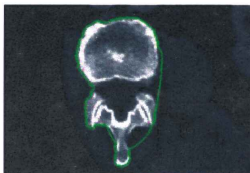
- Minimize the total mass under rigidity constraints.
- Shape optimization: maximize the stability (oscillations due to car traffic, pedestrian or wind) of hanging bridges.



Example 6: Image segmentation using active contours



Initialisation



1.



2.



3.



4.



Résultat

Example 6: Image segmentation using active contours

A contour is defined as a parametric function $M = f(s)$.

- The contours lays in the contrast areas of the image.

⇒ maximize

$$\int_s \|\nabla A(f(s))\|_2^2 ds$$

with respect to f .

- Find a smooth contour.

⇒ constraint that the smoothness function $R(f)$ is bounded.

- Formulation:

$$\min_f - \int_s \|\nabla A(f(s))\|_2^2 ds \quad \text{subject to } R(f) \leq A$$

Example 7: Image registration

$$\max_T S(I_t, T(I_s))$$

- T : geometrical transformation (translation, rotation, *etc.*).
- I_s : source image, I_t target image.
- S : measure of similarity.

Conclusion (part 1)

There is **no generic numerical** method to solve all kinds of problems!

The effectiveness of numerical optimization methods depends on:

- The differentiability of f .
- The possibility to compute efficiently its derivatives.
- The feasible set: discrete or continuous?
- The kinds of constraints (equalities / inequalities).
- The specific class of problem (linear programming, convex, *etc.*).
- The problem dimension.

Basic definitions

Local/global minimizers

$x^* \in \mathcal{D}$ is a **global minimizer** of f over \mathcal{D} if $\forall x \in \mathcal{D}, f(x) \geq f(x^*)$.

$x \in \mathcal{D}$ is a **local minimizer** of f over \mathcal{D} if there exists a neighborhood $\mathcal{V}(x)$ of x such $\forall y \in \mathcal{V}(x) \cap \mathcal{D}, f(y) \geq f(x)$.

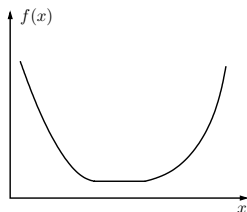
The **minimum** of f over \mathcal{D} is the value $\min_{x \in \mathcal{D}} f(x)$.

Unimodality / multimodality

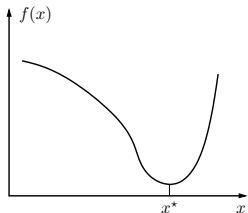
f is **unimodal** if f possesses a unique local minimizer. Otherwise, f is **multimodal**.

Remark.

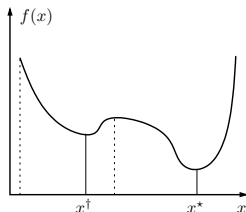
- **Strict** convexity \Rightarrow unimodality.
- Convexity \Rightarrow all local minimizers are global minimizers.



convex
but multimodal



nonconvex
but unimodal



nonconvex
and multimodal

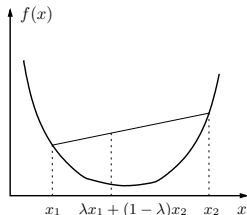
Convexity

f is **convex** if for any $x_1, x_2 \in \mathcal{D}$, and any $\lambda \in [0, 1]$,

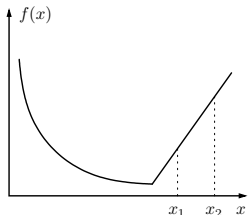
$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

f is **strictly convex** if for any $x_1, x_2 \in \mathcal{D}$, and any $\lambda \in]0, 1[$,

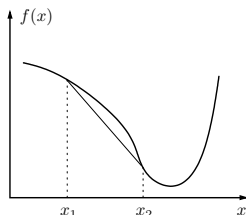
$$f(\lambda x_1 + (1 - \lambda)x_2) < \lambda f(x_1) + (1 - \lambda)f(x_2)$$



strictly convex



convex



nonconvex

Convexity: specific cases

Linear programming.

- $f(x) = c^T x = \sum_i c_i x_i$ is convex.

Quadratic programming.

- $f(x) = \|y - Ax\|_2^2 = \sum_i (y_i - a_i^T x)^2$ is convex.
- $f(x) = \frac{1}{2} x^T A x - b^T x = \frac{1}{2} \sum_i \sum_j a_{ij} x_j x_j - \sum_i b_i x_i$ is convex
when H is symmetric positive semidefinite.

Iterative algorithms

Iterative algorithms.

Solve the minimization with an iterative algorithm which generates a sequence of iterates x_0, x_1, \dots, x_n .

Convergence.

The sequence of iterates x_n converges to the solution x^* of the minimization problem.

- Global convergence: valid for all $x_0 \in \mathcal{D}$.
- Local convergence: valid for x_0 in a neighborhood of x^* .

Rate of convergence

- **Linear** convergence: $\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} \leq r$ for k sufficiently large.
- **Quadratic** convergence: $\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^2} \leq r$ for k sufficiently large.
- **Super-linear** convergence: $\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 0$.

Linear programming

Canonical form:

$$\text{Minimize } c^T x \quad \text{subject to } \begin{cases} Ax = b \\ x \geq \mathbf{0} \end{cases}$$

Standard form:

$$\text{Minimize } c^T x \quad \text{subject to } Ax \leq b$$

Remarks:

- Set $z = b - Ax \geq \mathbf{0}$. Standard form \Rightarrow Canonical form.
- No closed-form solution.

Resolution:

- The **simplex** algorithm.
- Integer linear programming (discrete case).

Linear least squares

Given $\mathbf{y} \in \mathbb{R}^m$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$, find $\boxed{\mathbf{x}^* \in \arg \min_x \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2}$

\mathbf{x}^* is a solution of

$$\begin{aligned}\nabla_x \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 &= \mathbf{0} \\ \Leftrightarrow 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{y}) &= \mathbf{0} \\ \Leftrightarrow \mathbf{A}^T\mathbf{A}\mathbf{x} &= \mathbf{A}^T\mathbf{y}\end{aligned}$$

If $m \geq n$ and \mathbf{A} is full column rank: $\mathbf{x}^* = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{y}$.

Otherwise:

- **infinity of solutions.**
- Generalized inverse: solution having the minimum ℓ_2 -norm:

Example. Find the best coefficients (a, b) of a linear model so as to reproduce at best observations y_i at times t_i .

Example. Data fitting using a polynomial of degree p .

Quadratic programming

Given $c \in \mathbb{R}^n$, $Q \in \mathbb{R}^{n \times n}$, $A \in \mathbb{R}^{p \times n}$,

Minimize $\frac{1}{2} x^T Q x + c^T x$ **subject to** $Ax \leq b$.

- p : number of constraints.
- j -th constraint: $a_j^T x \leq b_j$.
- Q is positive semi-definite \implies convex problem.
- Q is positive definite: strictly convex problem.

Nonlinear least squares

Given $\mathbf{y} \in \mathbb{R}^m$ and $\varphi : \mathbb{R}^n \mapsto \mathbb{R}^m$,

$$\textbf{Minimize } \|\mathbf{y} - \varphi(\mathbf{x})\|_2^2$$

- No closed form solution
- Specific algorithms taking the structure of the problem into account (Levenberg-Marquardt)

Example

Given a set of points (x_i, y_i) , $i = 1, \dots, N$, approximate the set of points using a parametric curve $y = \varphi(x; \theta)$.

Example: $\varphi(x; \alpha_1, \alpha_2, \tau_1, \tau_2) = \alpha_1 \exp(x/\tau_1) + \alpha_2 \exp(x/\tau_2)$

with $\theta = \{\alpha_1, \alpha_2, \tau_1, \tau_2\}$.

1. Formulate the problem as an optimization problem.
2. Call the appropriate Matlab function.

cf. Matlab simulations.

2. Unconstrained local optimization

Working assumptions

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ is \mathcal{C}^1 or \mathcal{C}^2 .

Definitions.

- Gradient of criterion $f(x)$: $\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \dots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n$

- Hessian of criterion $f(x)$:

$$H(x) = \nabla^2 f(x) = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \frac{\partial^2 f}{\partial x_i \partial x_j} & \dots \\ \dots & \dots & \dots \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Existence of a local minimizer

Theorem. Given \mathcal{U} a closed subset of \mathbb{R}^n , and a continuous function $f : \mathcal{U} \subset \mathbb{R}^n \rightarrow \mathbb{R}$.

- If \mathcal{U} is bounded,
- or if $\lim_{\substack{\|x\| \rightarrow \infty \\ x \in \mathcal{U}}} f(x) = +\infty$ [coercivity],

then f admits a local minimum on \mathcal{U} .

Line search strategy

An iterative algorithm generates a sequence x_0, x_1, \dots, x_n .

How to move from x_k to x_{k+1} ?

- Use the information about function f at x_k .
- Possibly use information from earlier iterates x_{k-1}, x_{k-2} , *etc.*

Line search.

- Choose a direction d_k .
- Solve ([approximately](#)) the 1D minimization problem:

$$\min_{t \geq 0} f(x_k + td_k)$$

- Set $x_{k+1} = x_k + td_k$

Feasible direction

d is a feasible direction at $x \in \mathcal{D}$ if

$$\exists A > 0, \forall t \in [0, A], x + td \in \mathcal{D}$$

Descent direction

- d is a descent direction of f at $x \in \mathcal{D}$ if

$$\exists A > 0, \forall t \in [0, A], f(x + td) \leq f(x)$$

- First order Taylor expansion:

$$f(x + td) = f(x) + td^T \nabla f(x) + o(t)$$

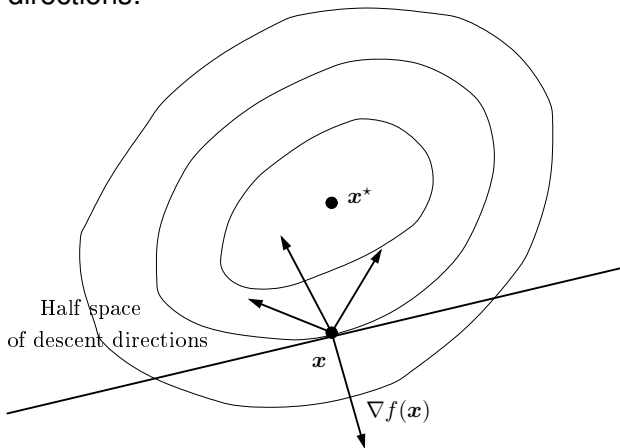
- d is a descent direction $\iff d^T \nabla f(x) < 0$

\implies Half space of descent directions.

Descent directions and level sets

Level sets: $\mathcal{L}_z = \{x \in \mathbb{R}^n, f(x) = z\}$.

Descent directions:



Derivatives

First order Taylor expansion:

$$f(x + td) = f(x) + td^T \nabla f(x) + o(t)$$

Second order Taylor expansion:

$$f(x + h) = f(x) + h^T \nabla f(x) + \frac{1}{2} h^T H(x) h + o(\|h\|^2)$$

$$f(x + td) = f(x) + td^T \nabla f(x) + \frac{t^2}{2} d^T H(x) d + o(t^2)$$

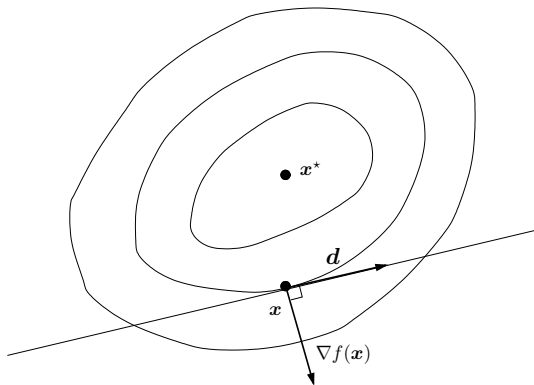
Directional derivative in the direction d :

$$\begin{aligned} f'(x; d) &= \lim_{t \rightarrow 0} \frac{f(x + td) - f(x)}{t} \\ &= d^T \nabla f(x) \end{aligned}$$

Derivatives

If d is tangent to the level set, then $f'(x; d) = d^T \nabla f(x) = 0$.

$\nabla f(x)$ is thus orthogonal to the level set.



Necessary condition for a local minimizer

Problem. Find $x^* \in \arg \min_{x \in \mathbb{R}^n} f(x)$

Note. $\max_{x \in \mathbb{R}^n} f(x) = - \min_{x \in \mathbb{R}^n} (-f(x))$

Theorem. If f admits a local minimum at x^* , then $\nabla f(x^*) = 0$.

Proof. Show that $f'(x^*; d) = d^T \nabla f(x^*) = 0 \ \forall d$.

Remark. The condition is necessary but not sufficient.

Sufficient condition for a local minimizer

Theorem. If $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\mathbf{H}(\mathbf{x}^*) > \mathbf{0}$, then f admits a local minimum at \mathbf{x}^* .

Positive definiteness. $\mathbf{H}(\mathbf{x}^*) > \mathbf{0} \Leftrightarrow \forall \mathbf{d} \neq \mathbf{0}, \mathbf{d}^T \mathbf{H}(\mathbf{x}^*) \mathbf{d} > 0$.

Proof. Show that $\forall \mathbf{d}, f(\mathbf{x}^* + t\mathbf{d}) > f(\mathbf{x}^*)$ for $|t|$ sufficiently small.

Descent algorithms

Choose x_0 .

For $k \geq 0$,

- Find a descent direction d_k
- Find a step $t_k > 0$.
- Compute $x_{k+1} = x_k + t_k d_k$.

Choice 1. Constant step $t_k = t, \forall k$.

Choice 2. **1D minimisation**: $t_k = \arg \min_{t>0} f(x_k + t d_k)$.

Descent algorithms

Descent: $\forall k, f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$.

Stopping conditions:

- $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) < \varepsilon_1$.
- $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 < \varepsilon_2$.
- $\|\nabla f(\mathbf{x}_k)\|_2 < \varepsilon_3$.
- $k = K_{\max}$.

Choice of the initial solution

- Unimodal criterion: any value of \mathbf{x}_0 should work (theoretically...)

Gradient based algorithms

The gradient algorithm.

- Choose x_0 .
- For $k \geq 0$, do $x_{k+1} \leftarrow x_k - t_k \nabla f(x_k)$.
 - \Rightarrow Constant step: $t_k = t, \forall k$.
 - \Rightarrow 1D minimization: $t_k \approx \arg \min_{t>0} f(x_k - t \nabla f(x_k))$.

The conjugate gradient algorithm.

- Choose x_0 .
- For $k \geq 0$, do $x_{k+1} \leftarrow x_k - t_k \nabla f(x_k) + \beta_k (x_k - x_{k-1})$.

Convergence rate is much faster!

Gradient algorithm

Remark.

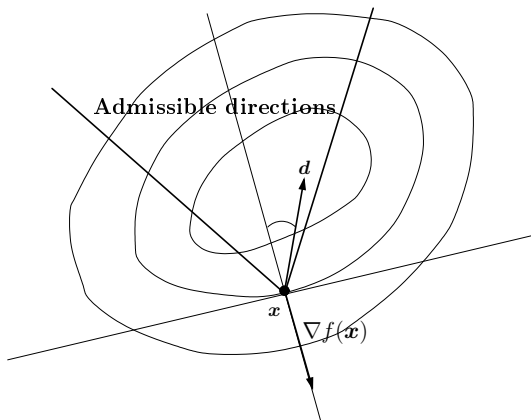
If the step is chosen as

$$t_k = \arg \min_{t>0} f(x_k - t \nabla f(x_k))$$

then for all k , $\nabla f(x_k)$ and $\nabla f(x_{k+1})$ are orthogonal (slow convergence, zigzag phenomenon).

⇒ Linear convergence.

Choice of descent direction



Impose that $|\text{Angle}(d_k, -\nabla f(x_k))| \leq \frac{\pi}{2} - \mu$.

Line minimization

How to compute $t_k \approx \arg \min_{t>0} f(x_k - t\nabla f(x_k))$?

- **Steepest descent**: exact minimization.
- **Constant step** $t_k = t$.
 - Risk of small steps.
 - Does not guarantee that $f(x_k + td_k) < f(x_k) \forall k$.
 \Rightarrow Reajust $t = t/2$.
- **Dichotomy**, golden section.
- **Quadratic approximation** of $\phi : t \rightarrow f(x_k + td_k)$.
- Quadratic or cubic approximation of ϕ from the knowledge of $\phi'(0)$.
- **Armijo and Wolfe's rules**: ensure that t is “acceptable” (sufficient decrease of f , sufficiently large step t).

Line minimization: example

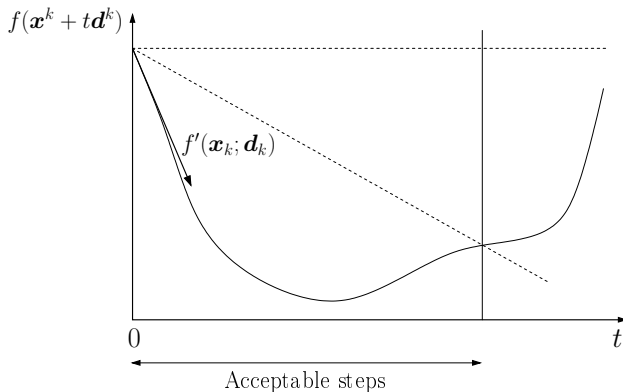
Quadratic interpolator.

Let $\phi(t) := f(\mathbf{x}_k + t\mathbf{d}_k)$.

Approximate $\phi(t) \approx q(t) = at^2 + bt + c$.

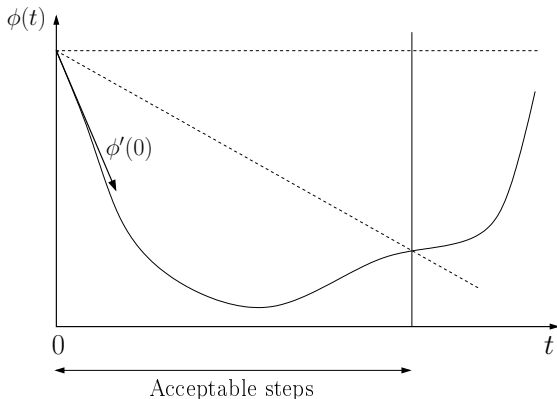
1. Compute $\phi(0) = f(\mathbf{x}_k)$.
2. Compute $\phi'(0) = \mathbf{d}_k^T \nabla f(\mathbf{x}_k)$.
3. Compute $\phi(t_0) = f(\mathbf{x}_k + t_0\mathbf{d}_k)$ for some t_0 .
4. Find the quadratic interpolator: compute a, b, c .
5. Minimize $q(t) \Rightarrow t_k$.

Step selection: Armijo's rule



Choose t such that $f(x_k + td_k) \leq f(x_k) + c_1 t f'(x_k; d_k)$
with $0 < c_1 < 1$.

Step selection: Armijo's rule

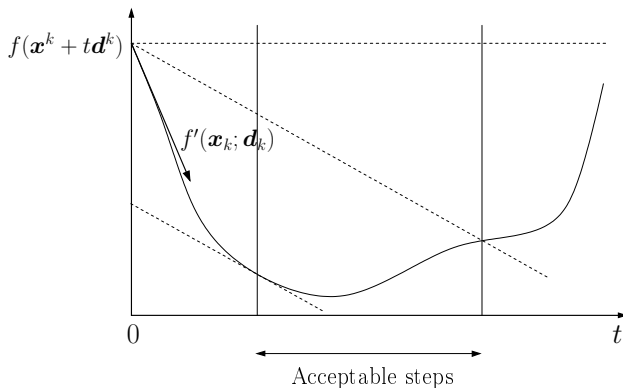


Let $\phi(t) := f(\mathbf{x}_k + t\mathbf{d}_k)$.

Armijo's condition rewrites $\phi(t) \leq \phi(0) + c_1 t \phi'(0)$

with $\phi'(0) = f'(\mathbf{x}_k; \mathbf{d}_k) = \nabla f(\mathbf{x}_k)^T \mathbf{d}_k$.

Step selection: Wolfe's rule



Choose t such that
$$\begin{cases} \phi(t) \leq \phi(0) + c_1 t \phi'(0) \\ \phi'(t) \geq c_2 \phi'(0) \text{ with } 0 < c_2 < c_1 < 1 \end{cases}$$

Step selection: Wolfe's rule

Choose t such that
$$\begin{cases} \phi(t) \leq \phi(0) + c_1 t \phi'(0) \\ \phi'(t) \geq c_2 \phi'(0) \text{ with } 0 < c_2 < c_1 \end{cases}$$

Rewriting:

Choose t such that

$$\begin{cases} f(x_k + t d_k) \leq f(x_k) + c_1 t \nabla f(x_k)^T d_k \\ \nabla f(x_k + t d_k)^T d_k \geq c_2 \nabla f(x_k)^T d_k \text{ with } 0 < c_2 < c_1 \end{cases}$$

Conjugate gradient algorithm

Conjugate gradient.

For $k \geq 0$,

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \beta_k \mathbf{d}_k \quad \text{with} \quad \mathbf{d}_k = -\nabla f(\mathbf{x}_k) + \beta_k \mathbf{d}_{k-1}.$$

- Fletcher-Reeves: $\beta_k^{\text{FR}} = \frac{\|\nabla f(\mathbf{x}_k)\|^2}{\|\nabla f(\mathbf{x}_{k-1})\|^2}$
- Polak-Ribière: $\beta_k^{\text{PR}} = \frac{[\nabla f(\mathbf{x}_k)]^T (\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}))}{\|\nabla f(\mathbf{x}_{k-1})\|^2}$

Subspace optimization (memory gradient).

$$\mathbf{d}_k = -\nabla f(\mathbf{x}_k) + \sum_{\ell=1}^m \omega_\ell \mathbf{d}_{k-\ell} \quad \text{with} \quad \omega_\ell > 0.$$

Use of second order derivatives

Idea.

1. Replace f by its quadratic approximation at point x_k :

$$f(x_k + d) \approx f(x_k) + d^T \nabla f(x_k) + \frac{1}{2} d^T H(x_k) d$$

2. Minimize the quadratic approximation: $d_k = -[H(x_k)]^{-1} g_k$

Works only if $H(x_k) > 0$.

Convex optimization

Newton algorithm.

- Choose x_0 .
- For $k \geq 0$, do $x_{k+1} \leftarrow x_k - t_k[H(x_k)]^{-1}\nabla f(x_k)$.

+ Fast (superlinear) convergence.

– Costly, not suited to large scale applications.

Adaptations

Newton direction: solve $\mathbf{H}(\mathbf{x}_k)\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$.

Truncated Newton direction: solve $\mathbf{H}(\mathbf{x}_k)\mathbf{d}_k \approx -\nabla f(\mathbf{x}_k)$:

find \mathbf{d}_k s.t. $\|\mathbf{H}(\mathbf{x}_k)\mathbf{d}_k + \nabla f(\mathbf{x}_k)\| \leq \eta_k \|\nabla f(\mathbf{x}_k)\|$ for $0 < \eta_k < 1$.

Modified Newton direction: when \mathbf{H}_k is not invertible, solve

$$[\mathbf{H}(\mathbf{x}_k) + \lambda_k \mathbf{I}_n]\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$$

for $\lambda_k > 0$.

Quasi-Newton methods

Idea: Approximate $H(x_k)$, $H(x_k)^{-1}$, or $H(x_k)^{-1}\nabla f(x_k)$ without computing second order derivatives.

Notations:

- $g_k := \nabla f(x_k)$.
- H_k : approximation of $H(x_k)$.
- B_k : approximation of $[H(x_k)]^{-1}$.

First order Taylor expansion:

$$\nabla f(x_{k+1} + h) = \nabla f(x_{k+1}) + H_{k+1}h + o(\|h\|)$$

\Rightarrow Choose H_{k+1} such that $g_{k+1} - g_k = H_{k+1}(x_{k+1} - x_k)$.

\Leftrightarrow Choose B_{k+1} such that $B_{k+1}(g_{k+1} - g_k) = x_{k+1} - x_k$.

Quasi-Newton algorithm

BFGS algorithm (Broyden, Fletcher, Goldfarb, Shanno):

$B_0 = I_n$ and:

$$B_{k+1} = B_k + \left(1 + \frac{p_k^T B_k p_k}{s_k^T p_k}\right) \frac{s_k s_k^T}{s_k^T p_k} - \frac{s_k p_k^T B_k + B_k p_k s_k^T}{s_k^T p_k}$$

with:

- $s_k := x_{k+1} - x_k$.
- $p_k := \nabla f(x_{k+1}) - \nabla f(x_k)$.

Notes.

- B_k is symmetric.
- For strictly convex f , BFGS preserves positive definiteness of B_k .

Quasi-Newton algorithm

DFP algorithm (Davidson, Fletcher, Powell):

$B_0 = I_n$ and:

$$B_{k+1} = B_k + \frac{s_k s_k^T}{s_k^T p_k} - \frac{B_k p_k p_k^T B_k}{p_k^T B_k p_k}$$

Quasi-Newton algorithm

Choose x_0 .

Set $B_0 = I_n$.

For $k = 0, 1, 2, \dots$,

- $d_k = -B_k g_k$.
- $t_k = \arg \min_{t>0} f(x_k + t d_k)$.
- $x_{k+1} = x_k + t_k d_k$.
- $B_{k+1} = B_k + \dots$ (BFGS rule).

End.

L-BFGS algorithm

Large scale problems, *e.g.*, image processing.

Avoid to store the Hessian matrix ($n(n + 1)/2$ coefficients).

Recursively update $B_k g_k$.

Matlab implementation: run fminunc

1. Write a Matlab function containing the variables x (here, called z) as the first input and at least one output (objective value f):

```
function [f,g] = myobj1(z,beta)

x=z(1);
y=z(2);

f = (y-sin(beta*x)-0.1*x*x)^2;
g(1) = -2*(beta*cos(beta*x)+0.2*x)*(y-sin(beta*x)-0.1*x*x);
g(2) = 2*(y-sin(beta*x)-0.1*x*x);
```

- Name of the file: same as the name of the function (here, `myobj1.m`).
- At least one output (f), gradient output g is optional.
- Recommendation: to speed up the optimization process, provide g as output when computation of g is possible.

Matlab implementation: run fminunc

2. **Run the Matlab solver** . This is another Matlab file, e.g., `runfminunc.m`:

```
clear all; % clear memory
close all; % close figures

beta = 2;

options = optimoptions('fminunc');
options = optimoptions(options,'Display','iter',...
    'MaxFunctionEvaluations',1000,'StepTolerance',1e-10,...
    'SpecifyObjectiveGradient',true,'CheckGradients',true);

% Initial solution
x0=[10;10];

% Compute initial objective value
fprintf('x0=(%.2f,%.2f), f=%f\n',x0(1),x0(2),...
    myobj1(x0,beta));
```

```
% Run optimization solver
```

```
xsol = fminunc(@(x)myobj1(x,beta),x0,options);
```

```
fprintf('Final solution=(%.2f,%.2f), f=%f\n', ...  
        xsol(1),xsol(2), myobj1(xsol,beta));
```

```
% 2D and 3D display
```

```
x=0:.1:6;
```

```
y=-1:.1:5;
```

```
% Pixel coordinates:
```

```
[xx,yy]=meshgrid(x,y);
```

```
zz=zeros(size(xx));
```

```
for i=1:length(xx(:)),
```

```
    zz(i)=myobj1([xx(i),yy(i)],beta);
```

```
end;
```



```
% 2D image
figure(1);
imagesc(x,y,zz);
colorbar;
axis xy
xlabel('x'); ylabel('y');
```

```
% 3D surface
figure(2);
mesh(x,y,zz);
colormap pink
shading flat
xlabel('x'); ylabel('y');
```

3. Least squares problems

Least squares problems

- Goal: $\min_{x \in \mathbb{R}^n} \left\{ f(x) = \frac{1}{2} \|F(x)\|_2^2 = \frac{1}{2} \sum_{j=1}^p F_j^2(x) \right\}$
- Residual vector $F : \mathbb{R}^n \mapsto \mathbb{R}^p$:

$$F(x) = \begin{bmatrix} F_1(x) \\ \vdots \\ F_p(x) \end{bmatrix}$$

- F linear with respect to x ?

Linear least squares vs nonlinear least squares.

Example 1

Predict the parameter of the model which best agrees with a time observation:

- y_j : observation at time t_j .
- p observations.
- $\phi(t; \mathbf{x}) = x_1 + tx_2 + t^2x_3$.

$$\Rightarrow \min_{\mathbf{x} \in \mathbb{R}^n} \sum_{j=1}^p [\phi(t_j; \mathbf{x}) - y_j]^2.$$

- Linear vs nonlinear least-squares?
 - $\mathbf{F}(\mathbf{x}) = \mathbf{J}\mathbf{x} - \mathbf{y}$?

Linear least-squares: $\min_{x \in \mathbb{R}^n} \left\{ f(x) = \frac{1}{2} \|F(x)\|_2^2 = \frac{1}{2} \|Jx - y\|_2^2 \right\}$

with:

$$F(x) = \begin{bmatrix} x_1 + t_1 x_2 + t_1^2 x_3 - y_1 \\ x_1 + t_2 x_2 + t_2^2 x_3 - y_2 \\ \vdots \\ x_1 + t_p x_2 + t_p^2 x_3 - y_p \end{bmatrix} \in \mathbb{R}^p$$

$$= \begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ \vdots & \vdots & \vdots \\ 1 & t_p & t_p^2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{bmatrix}$$

$$= Jx - y$$

- **Linear** least-squares problem: $\min_{x \in \mathbb{R}^n} \left\{ f(x) = \frac{1}{2} \|Jx - y\|_2^2 \right\}$
- Gradient: $\nabla f(x) = J^T(Jx - y)$.
- Solution:

$$\begin{aligned}\nabla f(x^*) = \mathbf{0} &\iff J^T(Jx^* - y) = \mathbf{0} \\ &\iff (J^T J)x^* = J^T y \quad (\text{normal equations}). \\ &\iff x^* = (J^T J)^{-1} J^T y.\end{aligned}$$

Example 2

Predict the parameter of the model which best agrees with a time observation:

- y_j : observation at time t_j .
- P observations.
- $\phi(t; \mathbf{x}) = x_1 + tx_2 + t^2x_3 + x_4e^{-x_5t}$.

$$\Rightarrow \min_{\mathbf{x} \in \mathbb{R}^n} \sum_{j=1}^p [\phi(t_j; \mathbf{x}) - y_j]^2.$$

- Nonlinear least-squares

Nonlinear least-squares: $\min_{x \in \mathbb{R}^n} \left\{ f(x) = \frac{1}{2} \|F(x)\|_2^2 \right\}$

with:

$$F(x_1, x_2, x_3, x_4, x_5) = \begin{bmatrix} F_1(x_1, x_2, x_3, x_4, x_5) \\ F_2(x_1, x_2, x_3, x_4, x_5) \\ \vdots \\ F_p(x_1, x_2, x_3, x_4, x_5) \end{bmatrix}$$
$$= \begin{bmatrix} x_1 + t_1 x_2 + t_1^2 x_3 + x_4 e^{-x_5 t_1} - y_1 \\ x_1 + t_2 x_2 + t_2^2 x_3 + x_4 e^{-x_5 t_2} - y_2 \\ \vdots \\ x_1 + t_p x_2 + t_p^2 x_3 + x_4 e^{-x_5 t_p} - y_p \end{bmatrix} \in \mathbb{R}^p$$

Nonlinear least-squares

- Objective function: $f(x) = \frac{1}{2} \sum_{j=1}^p F_j^2(x)$.
- Partial derivatives: $\frac{\partial f}{\partial x_i}(x) = \sum_{j=1}^p \frac{\partial F_j}{\partial x_i}(x) F_j(x)$
- Gradient:

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial F_1}{\partial x_1}(x) & \frac{\partial F_2}{\partial x_1}(x) & \dots & \frac{\partial F_p}{\partial x_1}(x) \\ \frac{\partial F_1}{\partial x_2}(x) & \frac{\partial F_2}{\partial x_2}(x) & \dots & \frac{\partial F_p}{\partial x_2}(x) \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial F_1}{\partial x_p}(x) & \frac{\partial F_2}{\partial x_p}(x) & \dots & \frac{\partial F_p}{\partial x_p}(x) \end{bmatrix} \begin{bmatrix} F_1(x) \\ F_2(x) \\ \vdots \\ F_p(x) \end{bmatrix}$$

Gradient and Jacobian

Jacobian $J(x)$ is a matrix of size $p \times n$ whose rows identify with $\nabla F_j(x)$:

$$J(x) = \begin{bmatrix} \frac{\partial F_1}{\partial x_1}(x) & \frac{\partial F_1}{\partial x_2}(x) & \dots & \frac{\partial F_1}{\partial x_n}(x) \\ \frac{\partial F_2}{\partial x_1}(x) & \frac{\partial F_2}{\partial x_2}(x) & \dots & \frac{\partial F_2}{\partial x_n}(x) \\ \vdots & \vdots & \vdots & \\ \frac{\partial F_p}{\partial x_1}(x) & \frac{\partial F_p}{\partial x_2}(x) & \dots & \frac{\partial F_p}{\partial x_n}(x) \end{bmatrix}$$

$$\nabla f(x) = [J(x)]^T F(x).$$

Second order derivatives

- Partial derivatives: $\frac{\partial f}{\partial x_i} = \sum_{k=1}^p \frac{\partial F_k}{\partial x_i} F_k$

- Second order derivatives:

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \sum_{k=1}^p \frac{\partial F_k}{\partial x_i} \frac{\partial F_k}{\partial x_j} + \sum_{k=1}^p F_k \frac{\partial^2 F_k}{\partial x_i \partial x_j}.$$

- Hessian of f : $H(x) = \nabla^2 f(x) = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \frac{\partial^2 f}{\partial x_i \partial x_j} & \dots \\ \dots & \dots & \dots \end{bmatrix} \in \mathbb{R}^{n \times n}$

$$\Rightarrow \boxed{H(x) = J(x)^T J(x) + \sum_{k=1}^p F_k H_k}$$

Gauss-Newton method

The gradient and Hessian of $f(x) = \frac{1}{2} \sum_{j=1}^p F_j^2(x)$ read:

$$\nabla f(x) = [J(x)]^T F(x)$$

$$H(x) = J(x)^T J(x) + \sum_{k=1}^p F_k H_k$$

1. Approximate $H \approx J(x)^T J(x)$.
2. Compute descent direction by solving

$$[J(x_k)^T J(x_k)] d_k = -J(x_k)^T F(x_k)$$

3. Line search:

$$\min_{t>0} f(x_k + t d_k)$$

Gauss-Newton method

Validity.

- When residuals are small ($F(x) \approx \mathbf{0}$).
- When $F_j(x)$ are nearly linear.

Advantages.

- Numerical cost is cheap (no calculation of second order derivatives).
- When $J(x_k)$ has full rank and $\nabla f(x) \neq \mathbf{0}$, d_k is a descent direction for f .

The Levenberg-Marquardt method

1. Approximate $H \approx J(x)^T J(x) + \lambda I_n$ for some $\lambda \geq 0$.
2. Compute descent direction by solving

$$[J(x_k)^T J(x_k) + \lambda_k I_n] d_k = -J(x_k)^T F(x_k)$$

3. Line search:

$$\min_{t > 0} f(x_k + t d_k)$$

Comments.

- Case where $J(x)$ is rank-deficient, or nearly so.
- $\lambda_k = 0$: Gauss-Newton method.
- $\lambda_k \rightarrow \infty$: the gradient algorithm (slower but numerically accurate).
- There are automatic strategies for tuning λ_k .

Matlab implementation: run lsqnonlin

1. Write a Matlab function containing the variables x as the first input and at least one output (residual vector $F(x)$):

```
function [F,J] = residuals(x,t,y)

% residual vector
F = y-x(1)*exp(-t/x(2)); % same size as y

% Jacobian matrix of size (length(t) x length(x))
% Column 1: derivative of  $F_i = y_i - x_1 \exp(-t_i/x_2)$  w.r.t.  $x_1$ 
% Column 2: derivative of  $F_i = y_i - x_1 \exp(-t_i/x_2)$  w.r.t.  $x_2$ 
J = [-exp(-t/x(2)), -x(1)/(x(2)*x(2))*t.*exp(-t./x(2))];
```

- Name of the file: residuals.m.
- At least one output (F), Jacobian output J is optional.
- Recommendation: to speed up the optimization process, provide J as output when computation of J is easy.

Matlab implementation: run lsqnonlin

2. **Run the Matlab solver.** This is another Matlab file, *e.g.*, `fit_expo.m`:

```
% load data t and y, defined as column vectors
figure(1), clf, plot(t,y,'o','linewidth',2); grid on;

% Initial parameter values
x_init = [10,10];

opt = optimoptions('lsqnonlin');
opt = optimoptions(opt,'Display','iter');

% If Jacobian is provided as output of function "residuals"
opt = optimoptions(opt,'SpecifyObjectiveGradient',true);

% for debug only, checks that the Jacobian computation
% is correct:
% opt = optimoptions(opt,'CheckGradients',true);
```



```
% Run optimization solver
```

```
x = lsqnonlin(@(x) (residuals(x,t,y)),x_init,[],[],opt);
```

```
disp('Solution:'); fprintf('%.2f ',x'); disp(' ');
```

```
% Plot prediction of data y using model  $x_1 \exp(-t/x_2)$ 
```

```
figure(1), hold on,
```

```
plot(t,x(1)*exp(-t/x(2)),'-r','linewidth',2);
```

```
% One could also plot the prediction of data y
```

```
% using the initial parameters x_init which is less accurate
```

```
% plot(t,x_init(1)*exp(-t/x_init(2)),'-g','linewidth',2);
```

4. Constrained local optimization

Constrained minimization

$$\text{Minimize } f(x) \quad \text{subject to} \quad \begin{cases} g_i(x) \leq 0, & i = 1, \dots, p \\ h_j(x) = 0, & j = 1, \dots, q. \end{cases}$$

Case of linear constraints:

- Each constraint reads $g_i(x) = a_i^T x - b_i = \sum_{k=1}^n a_i(k)x(k) - b_i$.
- Matrix form: $g(x) = Ax - b \leq \mathbf{0}_p$ with $A \in \mathbb{R}^{p \times n}$.
- Similarly, write $h(x) = Cx - d$.
- **Matlab:** A, b, C, d are provided as inputs of `fmincon`.
- Case of lower ($x_k \geq \ell_k$) and upper bounds ($x_k \leq u_k$): compute corresponding a_i and b_i .

Constrained minimization

$$\text{Minimize } f(x) \quad \text{subject to} \quad \begin{cases} g_i(x) \leq 0, & i = 1, \dots, p \\ h_j(x) = 0, & j = 1, \dots, q. \end{cases}$$

Case of **nonlinear constraints** (Matlab):

- Write a function

```
[g,h] = function mycon(x);
```

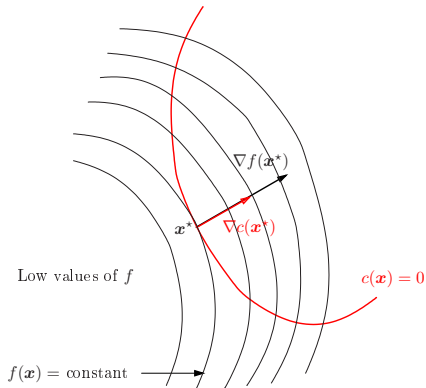
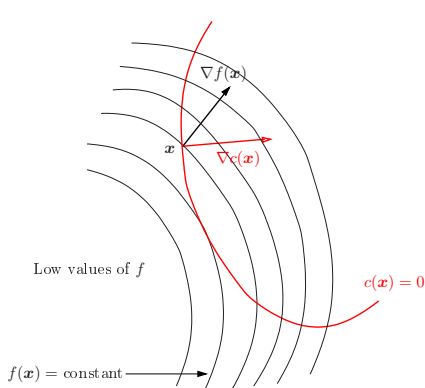
returning the values of $g_i(x)$ and $h_i(x)$ for all nonlinear constraints.
- Provide this function as input of `fmincon`.

Existence of a minimizer

We will assume that functions f , g_i and h_j are twice differentiable.

If the **feasible set** $\mathcal{D} = \{x \in \mathbb{R}^n, g(x) \leq \mathbf{0}_p, h(x) = \mathbf{0}_q\}$ is non-empty and bounded, then there exists a minimizer.

Case of 1 equality constraint $c(x) = 0$



At the optimal point x^* , $\nabla f(x^*) = \lambda \nabla c(x^*)$ for some λ .

Generalization for q equality constraints

Problem with **equality constraints only**:

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad \begin{cases} h_1(x) = 0 \\ h_2(x) = 0 \\ \vdots \\ h_q(x) = 0 \end{cases}$$

Define the **Lagrangian function** for $x \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}^q$:

$$\mathcal{L}(x, \lambda) = f(x) - \lambda_1 h_1(x) - \dots - \lambda_q h_q(x)$$

$$\lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_q \end{bmatrix} \text{ are Lagrange multipliers.}$$

First order optimality conditions

$$\text{Lagrangian: } \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \sum_{i=1}^q \lambda_i h_i(\mathbf{x})$$

If \mathbf{x}^* is a local minimizer of f , then there exists $\boldsymbol{\lambda}^*$ such that:

$$\begin{cases} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}_n \\ \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}_q \end{cases} \iff \begin{cases} \nabla f(\mathbf{x}^*) = \sum_{i=1}^q \lambda_i^* \nabla h_i(\mathbf{x}^*) \\ h_1(\mathbf{x}^*) = \dots = h_q(\mathbf{x}^*) = 0 \end{cases}$$

Example: $\min_{\mathbf{x} \in \mathbb{R}^2} (x_1^2 + x_2^2) \quad \text{s.t.} \quad x_1 + x_2 = 1$

1. Graphical representation.
2. Write the optimality conditions.
3. Solve the system.

Inequality constrained problems

Problem.

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad \begin{cases} g_1(x) \leq 0 \\ g_2(x) \leq 0 \\ \vdots \\ g_p(x) \leq 0 \end{cases}$$

Definition.

The i -th constraint is active at x if $g_i(x) = 0$.

The i -th constraint is inactive at x if $g_i(x) < 0$.

Lagrangian function

Problem with **inequality constraints only**:

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad \begin{cases} g_1(x) \geq 0 \\ g_2(x) \geq 0 \\ \vdots \\ g_p(x) \geq 0 \end{cases}$$

Define the Lagrangian function for $x \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}_+^p$:

$$\mathcal{L}(x, \lambda) = f(x) - \lambda_1 g_1(x) - \dots - \lambda_p g_p(x)$$

$$\lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_p \end{bmatrix} \geq \mathbf{0} \quad \text{are Lagrange multipliers.}$$

Lagrangian function (bis)

Modified Problem:

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad \begin{cases} g_1(x) \leq 0 \\ g_2(x) \leq 0 \\ \vdots \\ g_p(x) \leq 0 \end{cases}$$

Define the Lagrangian function for $x \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}_+^p$:

$$\mathcal{L}(x, \lambda) = f(x) + \lambda_1 g_1(x) + \dots + \lambda_p g_p(x)$$

$$\lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_p \end{bmatrix} \succeq \mathbf{0} \quad \text{are Lagrange multipliers.}$$

Karush-Kuhn-Tucker (KKT) conditions

If x^* is a local solution of

$$\min_x f(x) \quad \text{s.t.} \quad g(x) \geq \mathbf{0}_p$$

then there exists λ^* such that:

$$\left\{ \begin{array}{l} \nabla_{\mathbf{x}} \mathcal{L}(x^*, \lambda^*) = \mathbf{0}_n \iff \nabla f(x^*) = \sum_{i=1}^p \lambda_i^* \nabla g_i(x^*) \\ \forall i, g_i(x^*) \geq 0 \\ \forall i, \lambda_i^* \geq 0 \\ \forall i, \lambda_i^* g_i(x^*) = 0 \end{array} \right.$$

$$\left(\mathcal{L}(x, \lambda) = f(x) - \sum_{i=1}^p \lambda_i g_i(x) \right)$$

Karush-Kuhn-Tucker (KKT) conditions

If x^* is a local solution of

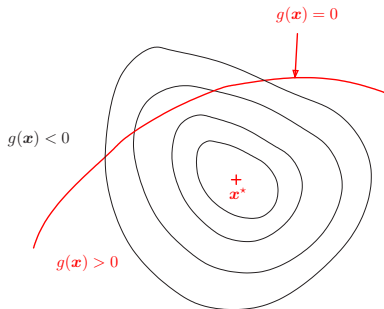
$$\min_x f(x) \quad \text{s.t.} \quad g(x) \leq \mathbf{0}_p$$

then there exists λ^* such that:

$$\left\{ \begin{array}{l} \nabla_{\mathbf{x}} \mathcal{L}(x^*, \lambda^*) = \mathbf{0}_n \iff \nabla f(x^*) = - \sum_{i=1}^p \lambda_i^* \nabla g_i(x^*) \\ \forall i, g_i(x^*) \leq 0 \\ \forall i, \lambda_i^* \geq 0 \\ \forall i, \lambda_i^* g_i(x^*) = 0 \end{array} \right.$$

$$\left(\mathcal{L}(x, \lambda) = f(x) + \sum_{i=1}^p \lambda_i g_i(x) \right)$$

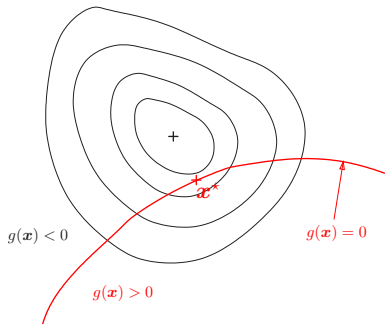
Case of 1 inequality constraint $g(x) \geq 0$



Inactive constraint: $g(x^*) > 0$

$$\lambda^* = 0$$

$$(\lambda^* g(x^*) = 0)$$

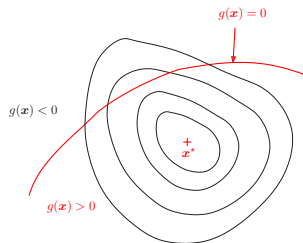


Active constraint: $g(x^*) = 0$

$$\lambda^* > 0$$

$$(\lambda^* g(x^*) = 0)$$

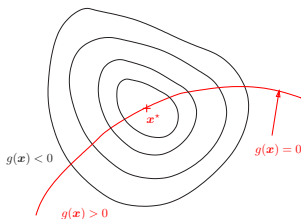
Inequality constraints



Inactive

$$\lambda^* = 0$$

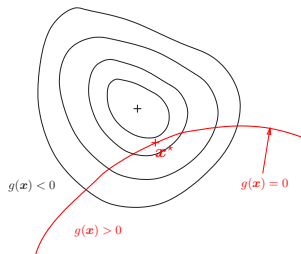
$$g(x^*) > 0$$



Weakly active

$$\lambda^* = 0$$

$$g(x^*) = 0$$



Active

$$\lambda^* > 0$$

$$g(x^*) = 0$$

Examples

Example 1: $\min_{x \in \mathbb{R}^2} x_1^2 + x_2^2 \quad \text{s.t.} \quad \begin{cases} x_1 \geq 0 \\ x_2 \geq 0 \\ x_1 + 2x_2 \geq 1. \end{cases}$

Example 2: [Nocedal & Wright, 2006, p. 329]

$$\min_{x \in \mathbb{R}^2} \left(x_1 - \frac{3}{2} \right)^2 + \left(x_2 - \frac{1}{8} \right)^4 \quad \text{s.t.} \quad \begin{cases} x_2 \leq 1 - x_1 \\ x_2 \geq -1 - x_1 \\ x_2 \geq x_1 - 1 \\ x_2 \leq x_1 + 1 \end{cases}$$

1. Graphical representation.
2. Write the optimality conditions.
3. Solve the system.

Necessary condition for a local minimizer

First order optimality condition.

If x^* is local minimizer, then x^* satisfies the KKT conditions.

Second order optimality condition.

- Let $\mathcal{A}(x)$ denote the set of active constraints at x
- Let $\mathcal{F}(x)$ denote the set of feasible directions at x .
- Let $\mathcal{C}(x, \lambda) = \{w \in \mathcal{F}(x) \mid [\nabla g_i(x)]^T w = 0 \forall i \in \mathcal{A}(x) \text{ with } \lambda_i > 0.\}$

If x^* satisfies the KKT conditions and

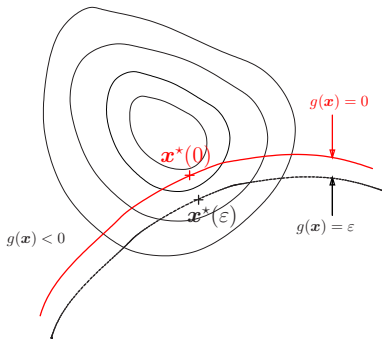
$$\forall w \in \mathcal{C}(x^*, \lambda^*) \setminus \{\mathbf{0}\}, w^T \nabla_{xx} \mathcal{L}(x^*, \lambda^*) w > 0$$

then x^* is a strict local minimizer of the inequality constrained problem.

Interpretation of Lagrange multipliers

For $\varepsilon > 0$, let $\mathbf{x}^*(\varepsilon) = \arg \min_{\mathbf{x}} f(\mathbf{x})$ s.t. $g(\mathbf{x}) \geq \varepsilon$.

Let $f^*(\varepsilon) = f(\mathbf{x}^*(\varepsilon))$. Then, $\lambda^* \|\nabla g(\mathbf{x}^*)\| = \frac{\partial f^*}{\partial \varepsilon}(\varepsilon = 0)$.



General case: inequality and equality constraints

Problem: $\min_{x \in \mathbb{R}^n} f(x)$ s.t. $\begin{cases} g_i(x) \geq 0, & i = 1, \dots, p \\ h_j(x) = 0, & j = 1, \dots, q. \end{cases}$

Define the Lagrangian function for $x \in \mathbb{R}^n$, $\lambda \in \mathbb{R}_+^p$, $\mu \in \mathbb{R}^q$:

$$\mathcal{L}(x, \lambda, \mu) = f(x) - \sum_{i=1}^p \lambda_i g_i(x) - \sum_{j=1}^q \mu_j h_j(x)$$

So, $\nabla_x \mathcal{L}(x, \lambda, \mu) = \nabla f(x) - \sum_{i=1}^p \lambda_i \nabla g_i(x) - \sum_{j=1}^q \mu_j \nabla h_j(x)$

KKT conditions

Problem: $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ s.t. $\begin{cases} g_i(\mathbf{x}) \geq 0, & i = 1, \dots, p \\ h_j(\mathbf{x}) = 0, & j = 1, \dots, q. \end{cases}$

If \mathbf{x}^* is a local minimizer, then there exists λ^* and μ^* such that:

$$\left\{ \begin{array}{l} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \lambda^*, \mu^*) = \mathbf{0}_n \\ \forall i, g_i(\mathbf{x}^*) \geq 0 \\ \forall j, h_j(\mathbf{x}^*) = 0 \\ \forall i, \lambda_i^* \geq 0 \\ \forall i, \lambda_i^* g_i(\mathbf{x}^*) = 0 \end{array} \right.$$

Example

$$\min_{\mathbf{x} \in \mathbb{R}^2} (4x_1^2 + 2x_2^2) \quad \text{s.t.} \quad \begin{cases} 3x_1 + x_2 = 8 \\ 2x_1 + 4x_2 \leq 15 \end{cases}$$

- $\mathcal{L}(x_1, x_2, \lambda, \mu) = 4x_1^2 + 2x_2^2 + \lambda(2x_1 + 4x_2 - 15) - \mu(3x_1 + x_2 - 8)$.
- $\frac{\partial \mathcal{L}}{\partial x_1}(x_1, x_2, \lambda, \mu) = 8x_1 + 2\lambda - 3\mu$.
- $\frac{\partial \mathcal{L}}{\partial x_2}(x_1, x_2, \lambda, \mu) = 4x_2 + 4\lambda - \mu$.

- KKT conditions:
$$\begin{cases} 8x_1 + 2\lambda - 3\mu = 0 \\ 4x_2 + 4\lambda - \mu = 0 \\ 3x_1 + x_2 = 8 \\ 2x_1 + 4x_2 \leq 15 \\ \lambda \geq 0 \\ \lambda(2x_1 + 4x_2 - 15) = 0 \end{cases}$$

$$\left\{ \begin{array}{rcl} 8x_1 + 2\lambda - 3\mu & = & 0 \\ 4x_2 + 4\lambda - \mu & = & 0 \\ 3x_1 + x_2 & = & 8 \\ 2x_1 + 4x_2 & \leq & 15 \\ \lambda & \geq & 0 \\ \lambda(2x_1 + 4x_2 - 15) & = & 0 \end{array} \right. \Rightarrow \left\{ \begin{array}{rcl} 8x_1 - 12x_2 & = & 10\lambda \\ 3x_1 + x_2 & = & 8 \\ 2x_1 + 4x_2 & \leq & 15 \\ \lambda & \geq & 0 \\ \lambda(2x_1 + 4x_2 - 15) & = & 0 \end{array} \right.$$

- Assume $2x_1 + 4x_2 = 15$. $3x_1 + x_2 = 8$ implies that $x_1 = \frac{17}{10}$, $x_2 = \frac{29}{10}$.
So, $\lambda = \frac{8x_1 - 12x_2}{10} = -\frac{212}{100} < 0$. **Contradiction!**
- So, $\lambda = 0$. The KKT conditions become:

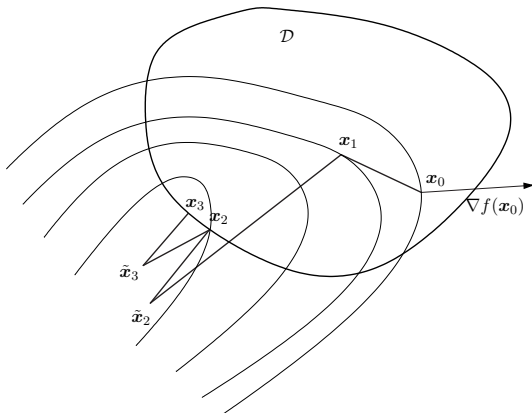
$$\left\{ \begin{array}{rcl} \lambda & = & 0 \\ \mu & = & \frac{8x_1}{3} \\ \mu & = & 4x_2 \\ 3x_1 + x_2 & = & 8 \\ 2x_1 + 4x_2 & \leq & 15 \end{array} \right. \Rightarrow \left\{ \begin{array}{rcl} \lambda & = & 0 \\ \mu & = & \frac{8x_1}{3} \\ x_1 & = & \frac{3}{2}x_2 \\ \frac{11}{2}x_2 & = & 8 \\ 2x_1 + 4x_2 & \leq & 15 \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{rcl} \lambda & = & 0 \\ \mu & = & \frac{8x_1}{3} \\ x_1 & = & \frac{24}{11} \\ x_2 & = & \frac{16}{11} \\ \text{Constraint } \frac{112}{11} & \leq & 15 \end{array} \right. \text{ (OK)}$$

Gradient projection algorithm

Principle.

1. Line minimization: $\tilde{x}_{k+1} = x_k + \alpha_k d_k$.
2. Projection onto convex set: $x_{k+1} = \text{proj}(\tilde{x}_{k+1}, \mathcal{D})$.



Other approaches

Interior point approaches.

Example: Replace $\min_{x \geq 0} f(x)$ by $\min_{x \in \mathbb{R}^n} \left\{ f(x) - \xi \sum_{i=1}^n \log x_i \right\}$.

Constrained minimization is replaced by the unconstrained minimization of the augmented criterion.

“Exterior” approaches, e.g., quadratic penalty, ADMM.

Example: Replace $\min_{x \geq 0} f(x)$ by

$$\min_{x, p \geq 0} \{ \mathcal{K}(x, p; \xi) = f(x) + \xi \|x - p\|^2 \} \quad \text{with } \xi \rightarrow \infty$$

Repeat:

- Minimize \mathcal{K} with respect to x (unconstrained).
- Minimize \mathcal{K} with respect to p (easy): $p_i = \max(x_i, 0)$.

Special case: linear programming

$$\min_{x \in \mathbb{R}^n} c^T x \quad \text{s.t.} \quad \begin{cases} Ax = b \\ x \geq \mathbf{0} \end{cases}$$

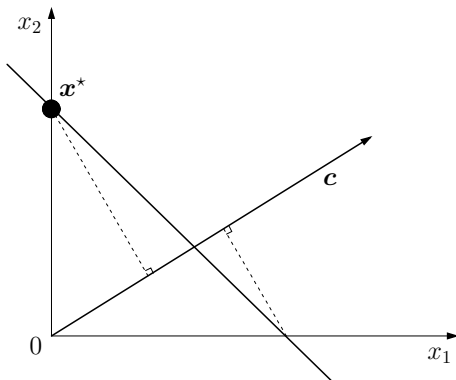
- $c \in \mathbb{R}^n$ with $c \geq \mathbf{0}_n$.
- $A \in \mathbb{R}^{m \times n}$ with $m \leq n$.

Example:

$$\min_{x \in \mathbb{R}^4} x_1 + 2x_2 + 3x_3 + 4x_4 \quad \text{s.t.} \quad \begin{cases} x_1 + x_2 + x_3 + x_4 = 1 \\ x_1 + x_3 - 3x_4 = \frac{1}{2} \\ x \geq \mathbf{0}. \end{cases}$$

- Rewrite the problem in matrix form.
- The solution x has two zero coordinates.

Interpretation in terms of orthogonal projections



The solution x^* has at least $(n - m)$ zero coordinates.

Case $n = 2, m = 1$: solution has at least one 0 coordinate.

Case $n = 3, m = 1$: solution has at least two 0 coordinates.

Case $n = 3, m = 2$: solution has at least one 0 coordinate.

Principle of the simplex algorithm

1. Decompose $A = [A_B, A_N]$ with A_B of size $m \times m$, and A_N of size $m \times (n - m)$.
2. Set $x^B = \begin{bmatrix} \mathbf{x}_B \\ \mathbf{0}_{n-m} \end{bmatrix}$ with $\mathbf{x}_B = A_B^{-1}b$ ⁽¹⁾.
3. x^B is a feasible point if $\mathbf{x}_B \geq \mathbf{0}_m$. The objective function is equal to $f(x^B) = c_B^T A_B^{-1}b$.
4. The **simplex algorithm** searches for the best decomposition $A = [A_B, A_N]$. The subset B is updated at each iteration in such a way that $f(x^B)$ is decreasing.

¹assuming that the rank of A equals m and that A_B is invertible

Special case: quadratic programming (QP)

$$\min_{x \in \mathbb{R}^n} Q(x) = \frac{1}{2} x^T H x + g^T x \quad \text{s.t.} \quad \begin{cases} a_i^T x = b_i, i \in \mathcal{E} \\ a_i^T x \geq b_i, i \in \mathcal{I} \end{cases}$$

- H is symmetric.
- H is symmetric positive semidefinite ($x^T H x \geq 0, \forall x$)
 \Rightarrow there exists a minimizer.
- H is symmetric positive definite ($x^T H x > 0, \forall x \neq \mathbf{0}$)
 \Rightarrow there exists a unique minimizer.

QP with equality constraints

$$\min_{x \in \mathbb{R}^n} Q(x) = \frac{1}{2} x^T H x + \ell^T x \quad \text{s.t. } Ax = b$$

1. Assume⁽²⁾ that $A \in \mathbb{R}^{m \times n}$ has rank m . Let $A = [A_1, A_2]$ with $A_1 \in \mathbb{R}^{m \times n}$ invertible, and let $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$. $Ax = b$ rereads $A_1 x_1 + A_2 x_2 = b$, so $x_1 = A_1^{-1}(b - A_2 x_2)$.
 \Rightarrow Solve unconstrained problem $\min_{x_2} Q(A_1^{-1}(b - A_2 x_2), x_2)$.
2. KKT conditions with $\mathcal{L}(x, \lambda) = \frac{1}{2} x^T H x + \ell^T x - \lambda^T (Ax - b)$:

Solve system

$$\begin{cases} Hx + \ell - A^T \lambda = \mathbf{0} \\ Ax = b \end{cases} \iff \begin{bmatrix} H & -A^T \\ A & \mathbf{0} \end{bmatrix} \begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} -\ell \\ b \end{bmatrix}$$

²without loss of generality

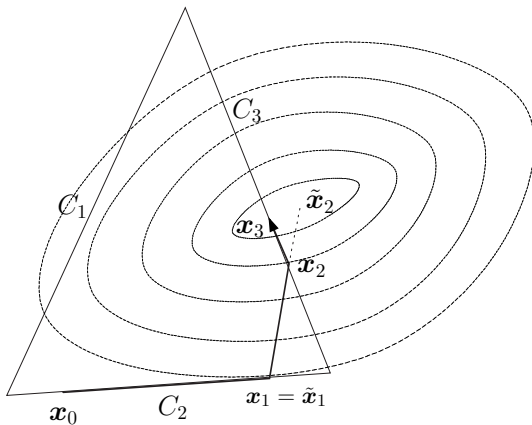
QP with inequality constraints

$$\min_{x \in \mathbb{R}^n} Q(x) = \frac{1}{2} x^T H x + \ell^T x \quad \text{s.t.} \quad Ax \geq b$$

Principle of the active-set algorithm.

1. Solve a sequence of QP problems with equality constraints.
2. Use a mechanism to add or remove active constraints.

The **active set** $\mathcal{A}(x) = \{i : a_i^T x = b_i\}$ is defined as the set of active constraints at x .



From a feasible point x_k :

1. Solve QP with equality constraints on $\mathcal{A}(x_k)$ only $\Rightarrow \tilde{x}_{k+1}$.
2. If \tilde{x}_{k+1} is feasible, set $x_{k+1} \leftarrow \tilde{x}_{k+1}$. Compute the Lagrange multipliers $\{\lambda_i, i \in \mathcal{A}(x_k)\}$. If some $\lambda_i < 0$, set $\mathcal{A}(x_{k+1}) \leftarrow \mathcal{A}(x_k) \setminus \{\arg \min_{i \in \mathcal{A}(x_k)} \lambda_i\}$ and go back to step 1. Else, STOP.
3. Otherwise, set $x_{k+1} \leftarrow x_k + \alpha_k(\tilde{x}_{k+1} - x_k)$ with the maximum possible value of $\alpha_k \in]0, 1[$. Compute $\mathcal{A}(x_{k+1})$ and go back to step 1.

Step 1

$$\text{QP: } \min_{x \in \mathbb{R}^n} Q(x) = \frac{1}{2} x^T H x + \ell^T x \quad \text{s.t. } Ax \geq b$$

QP with equality constraints: here, the inequality constraints are removed.

Define $A_{\mathcal{A}} = A(\mathcal{A}, :)$ and $\delta = x - x_k$ the descent direction.

$$\min_{x \in \mathbb{R}^n} Q(x) \quad \text{s.t. } A_{\mathcal{A}} x = b_{\mathcal{A}}$$

$$\Leftrightarrow \min_{\delta \in \mathbb{R}^n} Q(x_k + \delta) = Q(x_k) + Q(\delta) + x_k^T H^T \delta \quad \text{s.t. } A_{\mathcal{A}} \delta = \mathbf{0}$$

$$\Leftrightarrow \min_{\delta \in \mathbb{R}^n} \frac{1}{2} \delta^T H \delta + (H x_k + \ell)^T \delta \quad \text{s.t. } A_{\mathcal{A}} \delta = \mathbf{0}.$$

Step 2: check KKT conditions for QP

Lagrangian of $\min_{x \in \mathbb{R}^n} Q(x) = \frac{1}{2} x^T H x + \ell^T x \quad \text{s.t. } Ax \geq b :$

$$\mathcal{L}(x, \lambda) = \frac{1}{2} x^T H x + \ell^T x - \lambda^T A x + \lambda^T b$$

KKT conditions with $x = x_{k+1} \Rightarrow$

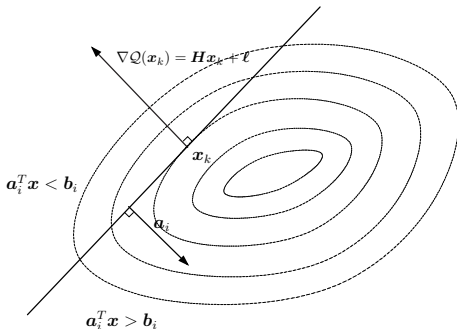
$$\begin{cases} Hx + \ell - A^T \lambda = \mathbf{0} \\ A_{\mathcal{A}} x = b_{\mathcal{A}} \\ \lambda_{\mathcal{A}} \geq \mathbf{0} \\ A_{\overline{\mathcal{A}}} x \geq b_{\overline{\mathcal{A}}} \\ \lambda_{\overline{\mathcal{A}}} = \mathbf{0} \end{cases}$$

$$A = \begin{bmatrix} A_{\mathcal{A}} \\ A_{\overline{\mathcal{A}}} \end{bmatrix} \text{ implies that } A^T \lambda = \begin{bmatrix} A_{\mathcal{A}}^T & A_{\overline{\mathcal{A}}}^T \end{bmatrix} \begin{bmatrix} \lambda_{\mathcal{A}} \\ \mathbf{0} \end{bmatrix} = A_{\mathcal{A}}^T \lambda_{\mathcal{A}}.$$

So, the KKT condition becomes

$$\lambda_{\mathcal{A}} = \left[A_{\mathcal{A}}^T \right]^{\dagger} (Hx + \ell) \geq \mathbf{0}$$

Case of an active-set with one constraint

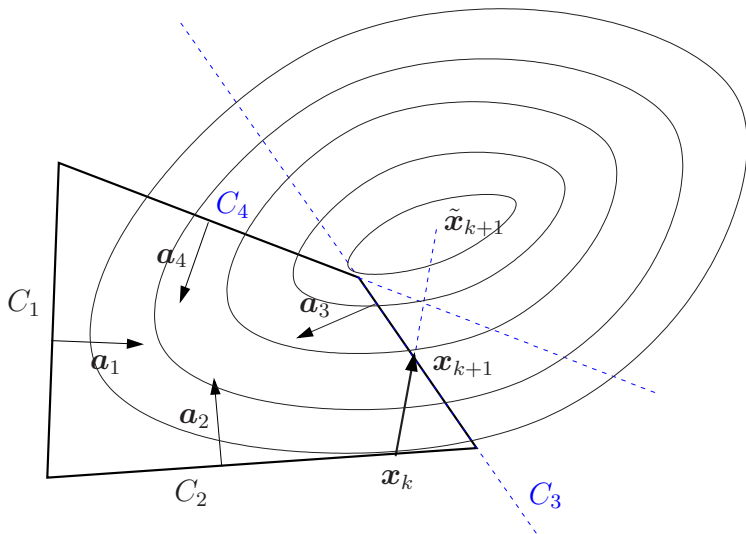


When $\mathcal{A} = \{i\}$ contains only one constraint, \mathbf{a}_i is colinear with $\nabla Q(\mathbf{x}_k) = \mathbf{H}\mathbf{x}_k + \boldsymbol{\ell}$, and

$$\lambda_i = [\mathbf{A}_{\mathcal{A}}^T]^\dagger (\mathbf{H}\mathbf{x}_k + \boldsymbol{\ell}) = [\mathbf{a}_i]^\dagger (\mathbf{H}\mathbf{x}_k + \boldsymbol{\ell}) = \frac{1}{\|\mathbf{a}_i\|^2} \mathbf{a}_i^T (\mathbf{H}\mathbf{x}_k + \boldsymbol{\ell}).$$

$\lambda_i < 0$ means that the i -th constraint may be removed

Step 3: illustration



Step 3: compute $x_k + \alpha_k(\tilde{x}_{k+1} - x_k)$

It is easy to check that $x_k + \alpha(\tilde{x}_{k+1} - x_k)$ is feasible when

$$\forall i \notin \mathcal{A}(x_k) \text{ such that } a_i^T(\tilde{x}_{k+1} - x_k) < 0, \alpha \leq \frac{b_i - a_i^T x_k}{a_i^T(\tilde{x}_{k+1} - x_k)}$$

So,

$$\alpha_k = \min_{\substack{i \notin \mathcal{A}(x_k) \\ a_i^T(\tilde{x}_{k+1} - x_k) < 0}} \frac{b_i - a_i^T x_k}{a_i^T(\tilde{x}_{k+1} - x_k)}$$

Sequential Quadratic Programming (SQP)

Solve a constrained problem $\min_{x \in \mathbb{R}^n} f(x)$ s.t. $\begin{cases} g_1(x) \geq 0 \\ \vdots \\ g_p(x) \geq 0 \end{cases}$

Repeat:

1. Replace f by its quadratic approximation around x_k :

$$f(x_k + \delta) \approx \mathcal{Q}_k(\delta) = f(x_k) + \nabla f(x_k)^T \delta + \frac{1}{2} \delta^T \nabla^2 \mathcal{L}(x_k, \lambda_k) \delta$$

2. Replace g_j by their linear approximation around x_k :

$$g_j(x_k + \delta) \approx g_j(x_k) + \nabla g_j(x_k)^T \delta$$

3. Solve the quadratic problem

$$\min_{x \in \mathbb{R}^n} \mathcal{Q}_k(\delta) \quad \text{s.t.} \quad \begin{cases} g_1(x_k) + \nabla g_1(x_k)^T \delta \geq 0 \\ \vdots \\ g_p(x_k) + \nabla g_p(x_k)^T \delta \geq 0 \end{cases}$$

in a “trust region”, i.e., a neighborhood of x_k .

5. Global optimization (recall definitions pp. 17-18)

Goal of global minimization

- Local algorithms: the numerical solution \hat{x} depends on the initial solution x_0 .
- Global optimization:
 - less sensitive to initialization.
 - escape from “poor valleys”.
 - often relies on the simulation of random variables: **stochastic algorithms**.

In practice

- First **try local solvers** with different initial solutions!
- Choose global optimization algorithms when the local optimization output strongly depends on the initial solution.
- **Never choose global optimization algorithms when the problem is unimodal** (especially for convex problems, including quadratic problems).
- Expectations: **find a local minimizer which is not too “poor”**.
Generally, global algorithms do not guarantee to find the global minimizer in a finite number of iterations.

Definitions

An algorithm is **deterministic** if the rules do not depend on random draws. For a given input x_0 , the algorithm will always produce the same iterates x_k , and the same output \hat{x} .

An algorithm is **stochastic** if random variables are generated during the iterations. Run the algorithm twice with the same input x_0 : the algorithm will not produce the same iterates x_k , nor the same output \hat{x} .

- Idea: **explore different valleys** where $f(x)$ is low as effectively as possible.
- Iterates are realizations of a random process.
- Define a probability density function.

Choice of a global optimization solver

Depends on:

- the size of the problem (number of variables).
- the number of local valleys.
- the time available for numerical computation.

By increasing level of complexity and computation time:

1. Random exploration with independent draws.
2. Interval optimization (deterministic).
3. Particle Swarm Optimization (PSO), Genetic Algorithm (GA).
4. Simulated Annealing.

Independent drawing of random variables

- Generate random points $x_k \in \mathcal{D}$.
- Evaluate $f(x_k)$.
- Select the point corresponding to the lowest value of $f(x_k)$.

$\hat{f} = +\infty.$

For $k = 1, \dots, K,$

Randomly generate $x_k \in \mathcal{D}$

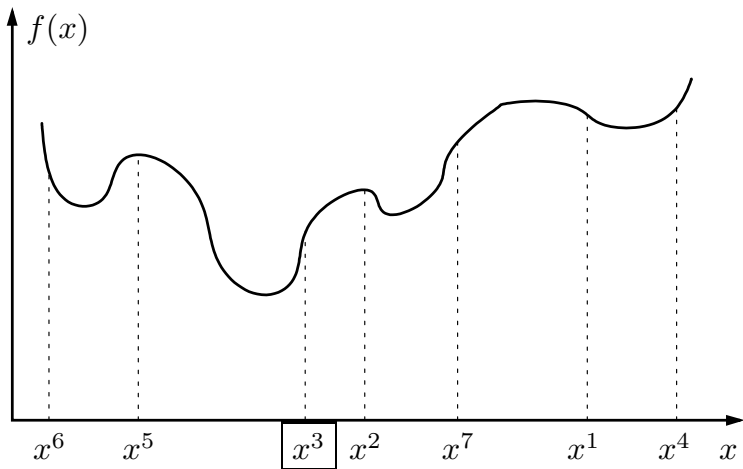
If $f(x_k) < \hat{f},$

$\hat{x} = x_k$ and $\hat{f} = f(x_k)$

End

End

Independent drawing of random variables



Hybrid approach: independent drawing + local optimization

$\hat{f} = +\infty.$

For $k = 1, \dots, K,$

Randomly generate $x_k \in \mathcal{D}$

Compute $y_k = \text{Local_Optim}(f; x_k)$

If $f(y_k) < \hat{f},$

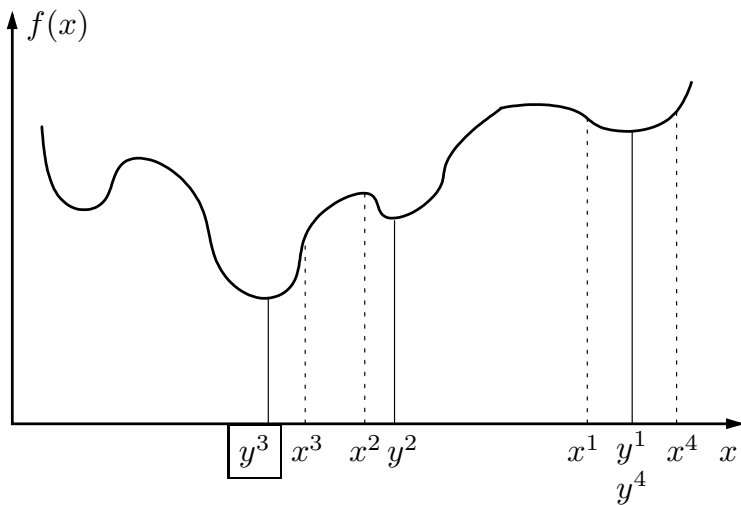
$\hat{x} = y_k$ and $\hat{f} = f(y_k)$

End

End

Note: y_k is the a local minimizer in the valley containing x_k .

Hybrid approach



Random drawing: which probability distribution?

$\hat{f} = +\infty.$

For $k = 1, \dots, K,$

Randomly generate $x_k \in \mathcal{D}$

If $f(x_k) < \hat{f},$

$\hat{x} = x_k$ and $\hat{f} = f(x_k)$

End

End For

- Effectiveness strongly depends on the choice of the probability distribution.
- Uniform distribution on \mathcal{D} ? (easy to simulate).
- More specific distribution if prior knowledge is available on the global minimizer x^* ?

Interval optimization

- A real number is represented by an interval, e.g., π is replaced by the interval $[3.14159, 3.14160]$.
- In 2D, an interval is a rectangle $[x_1^-, x_1^+] \times [x_2^-, x_2^+]$ with two lower bounds and two upper bounds.
- In n D, an interval is a cube $[x_1^-, x_1^+] \times [x_2^-, x_2^+] \times \cdots \times [x_n^-, x_n^+]$ of \mathbb{R}^n .
- Evaluate the objective function on an interval, that is, evaluate the image of an interval by function f . Examples:

$$[-2, 3] + [5, 7] = [3, 10]$$

$$[-3, 2] * [-3, 2] = [-6, 9]$$

$$[-3, 2]^2 = [0, 9]$$

$$\exp([-2, 3]) = [\exp(-2), \exp(3)]$$

$$\sin([0, 2\pi/3]) = [0, 1]$$

$$[-3, 4] - [-3, 4] = [-7, 7]$$

$$\ln([-2, -1]) = \emptyset$$

Evaluation of a function on an interval

- Decompose the function as a combination of simple functions, and evaluate each simple function.
- Example with $f(x) = x^2 + 2x + 4$:

$$\begin{aligned} f([-3, 4]) &= [-3, 4]^2 + 2 * [-3, 4] + [4, 4] \\ &= [0, 16] + [-6, 8] + [4, 4] = [-2, 28]. \end{aligned}$$

means that $\forall x \in [-3, 4], f(x) \in [-2, 28]$.

- Remark: the bounds -2 and 28 are not necessarily reached!
- Interval bisection: cut the interval in 2 sub-intervals, $[-3, 0.5]$ and $[0.5, 4]$ and evaluate each sub-interval.

Bisection: 2D example

We evaluate $f([x_1^-, x_1^+] \times [x_2^-, x_2^+]) = [-20, 200]$.

Then, the rectangle is split into 4 rectangles:

- $\mathcal{P}_{11} = [x_1^-, \frac{x_1^- + x_1^+}{2}] \times [x_2^-, \frac{x_2^- + x_2^+}{2}]$,
- $\mathcal{P}_{12} = [x_1^-, \frac{x_1^- + x_1^+}{2}] \times [\frac{x_2^- + x_2^+}{2}, x_2^+]$,
- $\mathcal{P}_{21} = [\frac{x_1^- + x_1^+}{2}, x_1^+] \times [x_2^-, \frac{x_2^- + x_2^+}{2}]$;
- $\mathcal{P}_{22} = [\frac{x_1^- + x_1^+}{2}, x_1^+] \times [\frac{x_2^- + x_2^+}{2}, x_2^+]$

and f is re-evaluated on each rectangle:

| | $[x_2^-, \frac{x_2^- + x_2^+}{2}]$ | $[\frac{x_2^- + x_2^+}{2}, x_2^+]$ |
|------------------------------------|------------------------------------|------------------------------------|
| $[x_1^-, \frac{x_1^- + x_1^+}{2}]$ | $f(\mathcal{P}_{11}) = [32, 100]$ | $f(\mathcal{P}_{12}) = [-14, 28]$ |
| $[\frac{x_1^- + x_1^+}{2}, x_1^+]$ | $f(\mathcal{P}_{21}) = [-10, 2]$ | $f(\mathcal{P}_{22}) = [2, 3]$ |

Bisection: 2D example

| | $\left[x_2^-, \frac{x_2^- + x_2^+}{2} \right]$ | $\left[\frac{x_2^- + x_2^+}{2}, x_2^+ \right]$ |
|---|---|---|
| $\left[x_1^-, \frac{x_1^- + x_1^+}{2} \right]$ | x | $f(\mathcal{P}_{12}) = [-14, 28]$ |
| $\left[\frac{x_1^- + x_1^+}{2}, x_1^+ \right]$ | $f(\mathcal{P}_{21}) = [-10, 2]$ | $f(\mathcal{P}_{22}) = [2, 3]$ |

Conclusion:

- \mathcal{P}_{11} cannot contain a global minimizer.
- \mathcal{P}_{12} , \mathcal{P}_{21} and \mathcal{P}_{22} are candidates to contain a global minimizer.
- Make bisection of \mathcal{P}_{12} , \mathcal{P}_{21} and \mathcal{P}_{22} and re-evaluate f on the new (smaller) intervals.

Interval optimization

Principle.

- Handle of list of candidate intervals $\{\mathcal{P}_i\}$.
- Evaluate $[f_i^-, f_i^+] = f(\mathcal{P}_i)$.
- Update $\bar{f} = \min_i f_i^+$.
- Simplification of the list: remove \mathcal{P}_i from the list when $f_i^- > \bar{f}$.

Branch and bound philosophy.

- All intervals are tested (bisection + evaluation of f), but some which cannot lead to the solution are deleted.
- The algorithm works for problems of dimension $n \leq 100$.
- The algorithm output is a list of candidate intervals and an upper bound \bar{f} of $\min f(x)$.
- The global minimizers are guaranteed to be in one of the intervals in the list.

Evolutionary algorithms

- Particle Swarm Optimization.
- Genetic Algorithm.

Principle.

Evolutionary algorithms make an exploration of the feasible domain \mathcal{D} using **several variables** in parallel.

Local optimization algorithm.

- 1 variable x .
- Iteration k : $x_{k+1} = x_k + t_k d_k$.

Evolutionary algorithms.

- N variables $x^1, \dots, x^N \in \mathbb{R}^n$.
- Iteration k : **jointly** update $\{x_{k+1}^1, \dots, x_{k+1}^N\}$ from $\{x_k^1, \dots, x_k^N\}$.

Evolutionary algorithms

Find a strategy with several variables (individuals, agents, particles) to explore “interesting valleys”. The number of variables N is fixed.

The update of the x_k^i depends on:

1. a local descent move (knowledge of $\nabla f(x_k^i)$).
2. the introduction of random moves to avoid convergence towards a local minimizer.
3. the objective values for other variables x_k^j ($j \neq i$).

Heuristic algorithms, no guarantee of convergence towards a global minimizer.

Simple to implement, several parameters to tune, can work for large dimension n .

Genetic algorithms

Artificial intelligence: simulate the evolution of a population of N individuals according to Darwin's theory:

1. **Natural selection**: the most adapted individuals tend to live longer and to reproduce more easily.
2. **Evolution**: some novel species randomly appear.

Sketch of algorithm:

Set $k \leftarrow 0$.

Initialize $x_0^1, x_0^2, \dots, x_0^N \in \mathbb{R}^n$

For $k = 0, \dots, K$,

 Select $N/2$ variables (the others are removed);

 Reproduction (mutations, crossbreed);

 Evaluate children;

 Replace some parents by their children.

End For

Genetic algorithms

Initialization.

x_0^i : uniform distribution on \mathcal{D} ?

Selection:

1. *N/2-elitism*: keep the $N/2$ best variables (corresponding to the lowest values of $f(x)$).
2. *Tournament selection*:

Repeat:

- Choose q variables randomly;
- Select the best as parent.

until the number of parents is equal to $N/2$.

Remark.

- *N/2-elitism* may induce premature convergence towards a local minimizer.
- Tournament selection will tend to explore unlikely valleys.

Genetic algorithms

Crossbreed between variables x_i^k and x_j^k :

Generate a new position

$$x = \alpha_{ij}x_k^i + (1 - \alpha_{ij})x_k^j$$

where α_{ij} is uniformly drawn on the line joining x_i^k and x_j^k . For instance, choose $\alpha_{ij} \in [-\varepsilon, 1 + \varepsilon]$.

Mutation of x_k^i :

$$x = x_k^i + \varepsilon_i \text{ with } \varepsilon_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

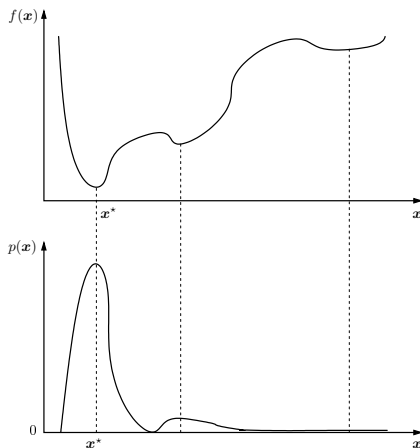
Overall parameters.

- N , σ , q (tournament selection), initial conditions.
- Tradeoff exploration (low q , large σ , large N) vs local convergence.

Monte Carlo based methods (simulated annealing...)

Minimize $f(x) \iff$ Maximize $p(x) = \frac{1}{Z} \exp(-f(x))$.

$Z = \int_{\mathcal{D}} \exp(-f(x)) dx$ is chosen to have a probability distribution.








Monte Carlo based methods

Minimize $f(x)$ \iff Maximize $p(x) = \frac{1}{Z} \exp(-f(x))$.

- Draw samples x_k of the probability distribution $p(x)$.
- For large dimension problems, and when Z is unknown, we can make use of Markov chains.
- MCMC methods (*Markov chain Monte Carlo*) create a long Markov chain $\{\dots, \mathbf{X}^k, \dots\}$ whose samples are asymptotically distributed according to pdf $p(x)$.

Bibliography

-  J. Nocedal and J. S. Wright,
Numerical optimization, 2nd edition, Springer Verlag, New York, Jul. 2006.
-  R. Fletcher,
Practical Methods of Optimization, 2nd edition, Wiley, 2000.
-  A. Björck,
Numerical Methods for Least Squares Problems, *Society for Industrial and Applied Mathematics*, Philadelphia, Apr. 1996.
-  R. Horst, P. M. Pardalos, and V. T. Nguyen,
Introduction to Global Optimization, 2nd edition, Springer, Dec. 2000.
-  P. M. Pardalos, D.-Z. Du, and R. L. Graham,
Handbook of Combinatorial Optimization, Springer Verlag, New York, 2013.

Acknowledgements

Many thanks to:

- Pierre Riedinger (University of Lorraine).
- Said Moussaoui (Ecole Centrale Nantes).
- Jérôme Idier (CNRS).