

Non-Linear Optimization Methods - Part I

Optimization Methods in Management Science

Master in Management

HEC Lausanne

Dr. Rodrigue Ouevray

Fall 2019 Semester

Non-Linear Optimization Methods

Unconstrained optimization:

- Quadratic programming
 - ▶ Direct resolution
 - ▶ Conjugate gradients

Differentiable Optimization

- When the objective function is **smooth**, numerical methods using information about **derivatives** (first order and possibly second order) perform better than derivative-free algorithms
- In the coming slides, we will focus on **quadratic programming**. For a quadratic function, it is **easy** to compute its gradient and its hessian
- Quadratic programming has many applications. Two applications will be discussed in this course: **Portfolio Optimization** and **Support Vector Machine**

Quadratic Programming

- Suppose we want to solve the following problem Q :

$$\min_{\mathbf{x} \in \mathbb{R}^n} \quad \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{g}^T \mathbf{x} + c,$$

where \mathbf{Q} is a **symmetric positive definite** matrix

- They are many different methods to solve this problem. We present here a direct resolution of it
- The solution of this system is simply given by (first order conditions):

$$\mathbf{Q} \mathbf{x} + \mathbf{g} = \mathbf{0} \iff \mathbf{Q} \mathbf{x} = -\mathbf{g} \iff \mathbf{x} = -\mathbf{Q}^{-1} \mathbf{g}$$

- **Solving an unconstrained quadratic problem defined by a symmetric positive definite matrix is equivalent to solving a linear system !**

Direct Resolution of a Quadratic Program (1)

- A direct resolution is typically based on the **Cholesky decomposition** of matrix \mathbf{Q}
- A **Cholesky decomposition** of a **symmetric positive definite** matrix \mathbf{Q} is defined as $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$, where \mathbf{L} is a **lower triangular** matrix
- A square matrix is called **lower triangular** if all the entries above the main diagonal are zero
- **This decomposition only works when the matrix is positive definite !**

Cholesky Decomposition : Example

- We would like to determine the Cholesky decomposition of \mathbf{Q} :

$$\begin{aligned}\mathbf{Q} &= \begin{pmatrix} 1 & 0 & 3 \\ 0 & 4 & 2 \\ 3 & 2 & 11 \end{pmatrix} = \begin{pmatrix} L_{11} & 0 & 0 \\ L_{21} & L_{22} & 0 \\ L_{31} & L_{32} & L_{33} \end{pmatrix} \begin{pmatrix} L_{11} & L_{21} & L_{31} \\ 0 & L_{22} & L_{32} \\ 0 & 0 & L_{33} \end{pmatrix} \\ &= \begin{pmatrix} L_{11}^2 & * & * \\ L_{11}L_{21} & L_{21}^2 + L_{22}^2 & * \\ L_{11}L_{31} & L_{21}L_{31} + L_{22}L_{32} & L_{31}^2 + L_{32}^2 + L_{33}^2 \end{pmatrix}\end{aligned}$$

- By identification, it is easy to determine the values of the L_{ij} :
 - ▶ row 1, column 1: $L_{11} = 1$
 - ▶ row 2, column 1: $L_{21} = 0$
 - ▶ row 2, column 2: $L_{22} = 2$
 - ▶ ...

Cholesky Decomposition : Example (Cont'd)

We finally get the following decomposition :

$$\mathbf{Q} = \begin{pmatrix} 1 & 0 & 3 \\ 0 & 4 & 2 \\ 3 & 2 & 11 \end{pmatrix} = \begin{pmatrix} \mathbf{1} & 0 & 0 \\ \mathbf{0} & \mathbf{2} & 0 \\ \mathbf{3} & \mathbf{1} & \mathbf{1} \end{pmatrix} \begin{pmatrix} \mathbf{1} & \mathbf{0} & \mathbf{3} \\ 0 & \mathbf{2} & \mathbf{1} \\ 0 & 0 & \mathbf{1} \end{pmatrix}$$

Direct Resolution of a Quadratic Program (2)

A **direct resolution** is based on the following steps:

- (1) Perform a Cholesky decomposition of $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$, where \mathbf{L} is a lower triangular matrix
- (2) Compute \mathbf{y} the solution of $\mathbf{L}\mathbf{y} = -\mathbf{g}$
- (3) Compute \mathbf{x} the solution of $\mathbf{L}^T\mathbf{x} = \mathbf{y}$

Then \mathbf{x} is the **unique** minimizer of the problem

Example

- We consider the following problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \quad \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{g}^T \mathbf{x}$$

where \mathbf{Q} and \mathbf{g} are given by:

$$\mathbf{Q} = \begin{pmatrix} 4 & 12 & -16 \\ 12 & 37 & -43 \\ -16 & -43 & 98 \end{pmatrix}, \quad \text{and} \quad \mathbf{g} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

- One can check that \mathbf{Q} is a symmetric positive definite matrix
- The Cholesky decomposition of \mathbf{Q} is

$$\mathbf{Q} = \begin{pmatrix} 4 & 12 & -16 \\ 12 & 37 & -43 \\ -16 & -43 & 98 \end{pmatrix} = \begin{pmatrix} 2 & 0 & 0 \\ 6 & 1 & 0 \\ -8 & 5 & 3 \end{pmatrix} \begin{pmatrix} 2 & 6 & -8 \\ 0 & 1 & 5 \\ 0 & 0 & 3 \end{pmatrix}$$

Example (Cont'd)

First, we need to solve :

$$\begin{pmatrix} 2 & 0 & 0 \\ 6 & 1 & 0 \\ -8 & 5 & 3 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = - \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \iff \begin{cases} 2y_1 & = & -1 \\ 6y_1 + y_2 & = & -1 \\ -8y_1 + 5y_2 + 3y_3 & = & -1 \end{cases}$$

We get that $y_1 = -0.5$, $y_2 = 2$, and $y_3 = -5$

Example (Cont'd)

Then, we need to find \mathbf{x} such that

$$\begin{pmatrix} 2 & 6 & -8 \\ 0 & 1 & 5 \\ 0 & 0 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} -0.5 \\ 2 \\ -5 \end{pmatrix} \iff \begin{cases} 2x_1 + 6x_2 - 8x_3 = -0.5 \\ x_2 + 5x_3 = 2 \\ 3x_3 = -5 \end{cases}$$

We conclude that the unique minimizer of this problem is given by:

$$x_1 = -\frac{455}{12}, x_2 = \frac{31}{3}, x_3 = -\frac{5}{3}$$

Quadratic Programming: The Conjugate Gradient Method

- Another approach to solve an **unconstrained** quadratic problem $f(\mathbf{x})$ when the matrix is **positive definite** is based on an **iterative** method called **conjugate gradient**
- It generates a **sequence of iterates** converging to the **optimal** solution of the problem
- The algorithm **stops** when a **stopping criterion** is met. It can be a **maximum number** of iterations and/or a criterion based on some necessary conditions to get an optimum (typically $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$, where ϵ is a small positive number)

The Conjugate Gradient Method

- Let \mathbf{Q} be a $n \times n$ **positive definite** matrix. We say that the non-null vectors \mathbf{d}_k are **Q-conjugate** if $\mathbf{d}_i^T \mathbf{Q} \mathbf{d}_j = 0, \forall i, j$ such that $i \neq j$
- The idea is to generate a sequence of iterates \mathbf{x}_k converging to the optimal solution:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k, \quad k = 1, \dots, n,$$

with $\alpha_k = \operatorname{argmin}_{\alpha} f(\mathbf{x}_k + \alpha \mathbf{d}_k)$

- $\alpha_k (\in \mathbb{R})$ is the **optimal** scalar that minimizes the objective function in the direction given by \mathbf{d}_k . Such α_k is determined by the following formula:

$$\alpha_k = -\frac{\mathbf{d}_k^T (\mathbf{Q} \mathbf{x}_k + \mathbf{g})}{\mathbf{d}_k^T \mathbf{Q} \mathbf{d}_k}$$

The Conjugate Gradient Method (Cont'd)

The directions \mathbf{d}_i are computed based on the **Gram-Schmidt process** applied to $-\nabla f(\mathbf{x}_1), \dots, -\nabla f(\mathbf{x}_i)$:

$$\mathbf{d}_i = -\nabla f(\mathbf{x}_i) + \sum_{k=1}^{i-1} \frac{\mathbf{d}_k^T \mathbf{Q} \nabla f(\mathbf{x}_i)}{\mathbf{d}_k^T \mathbf{Q} \mathbf{d}_k} \mathbf{d}_k$$

Indeed, one can directly compute \mathbf{d}_i from \mathbf{d}_{i-1} with the following formula :

$$\mathbf{d}_i = -\mathbf{Q}\mathbf{x}_i - \mathbf{g} + \beta_i \mathbf{d}_{i-1},$$

where

$$\beta_i = \frac{\nabla f(\mathbf{x}_i)^T \nabla f(\mathbf{x}_i)}{\nabla f(\mathbf{x}_{i-1})^T \nabla f(\mathbf{x}_{i-1})} = \frac{(\mathbf{Q}\mathbf{x}_i + \mathbf{g})^T (\mathbf{Q}\mathbf{x}_i + \mathbf{g})}{(\mathbf{Q}\mathbf{x}_{i-1} + \mathbf{g})^T (\mathbf{Q}\mathbf{x}_{i-1} + \mathbf{g})}$$

The CG Algorithm (1)

We consider an **unconstrained quadratic minimization** problem given by a **symmetric positive definite** matrix \mathbf{Q} and a vector \mathbf{g}

- **Input:** a first approximation \mathbf{x}_1 of the solution, a convergence parameter ϵ , and a parameter k counting the number of iterations
- **Output:** the (approximated) solution \mathbf{x}^* of the problem
- **Stopping criteria:**
 - ▶ If $k < n$ and $\|\nabla f(\mathbf{x}_{k+1})\| < \epsilon$, then $\mathbf{x}^* = \mathbf{x}_{k+1}$
 - ▶ If $k = n$, then $\mathbf{x}^* = \mathbf{x}_{n+1}$

The CG Algorithm (2)

- **Initialization:** $k = 1, \mathbf{d}_1 = -\mathbf{Q}\mathbf{x}_1 - \mathbf{g}$
- **Iterations:**
 - ▶ Compute the step α_k :

$$\alpha_k = -\frac{\mathbf{d}_k^T (\mathbf{Q}\mathbf{x}_k + \mathbf{g})}{\mathbf{d}_k^T \mathbf{Q}\mathbf{d}_k}$$

- ▶ Compute the next iterate $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$
- ▶ Check the stopping criterion. If it is satisfied, then stop: \mathbf{x}_{k+1} is the (approximated) solution of the problem. If not, then continue the current iteration
- ▶ Compute β_{k+1} :

$$\beta_{k+1} = \frac{\nabla f(\mathbf{x}_{k+1})^T \nabla f(\mathbf{x}_{k+1})}{\nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k)} = \frac{(\mathbf{Q}\mathbf{x}_{k+1} + \mathbf{g})^T (\mathbf{Q}\mathbf{x}_{k+1} + \mathbf{g})}{(\mathbf{Q}\mathbf{x}_k + \mathbf{g})^T (\mathbf{Q}\mathbf{x}_k + \mathbf{g})}$$

- ▶ Compute the new direction $\mathbf{d}_{k+1} = -\mathbf{Q}\mathbf{x}_{k+1} - \mathbf{g} + \beta_{k+1} \mathbf{d}_k$
- ▶ $k = k + 1$

Convergence of the CG Algorithm

- It **converges** in **at most** n iterations !
- For a quadratic problem defined by a symmetric positive definite matrix, a **sufficient** condition to have a minimum is that the gradient at that point should be **null**
- If the gradient is sufficiently close to **0** (the null vector), then we consider that we have found the optimum
- Concretely, if the **norm** of gradient evaluated at the current iterate is sufficiently small ($\|\nabla f(\mathbf{x}_{k+1})\| < \epsilon$), then we consider that **the algorithm has converged** to the optimum

Example

We consider the following problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \quad \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{g}^T \mathbf{x}$$

where \mathbf{Q} and \mathbf{g} are given by:

$$\mathbf{Q} = \begin{pmatrix} 4 & 1 \\ 1 & 3 \end{pmatrix}, \quad \mathbf{g} = \begin{pmatrix} -1 \\ -2 \end{pmatrix}$$

Example (Cont'd)

- Solving this problem is equivalent to solving the following system:

$$\begin{pmatrix} 4 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

- A direct resolution (with Gauss elimination or with the Cholesky decomposition) gives the following solution $\mathbf{x}^* \approx (0.09 \ 0.64)^T$
- Rather than solving directly this problem, let's try to compute the steps generated by the CG algorithm
- We assume that the starting point (chosen arbitrarily) for the CG algorithm is $\mathbf{x}_1 = (2 \ 1)^T$

Example: Iteration 1

- Initialization: $k = 1$, $\epsilon = 10^{-6}$,

$$\mathbf{d}_1 = -\mathbf{Q}\mathbf{x}_1 - \mathbf{g} = -\begin{pmatrix} 4 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} - \begin{pmatrix} -1 \\ -2 \end{pmatrix} = \begin{pmatrix} -8 \\ -3 \end{pmatrix}$$

- Compute the step α_1 :

$$\alpha_1 = -\frac{\mathbf{d}_1^T(\mathbf{Q}\mathbf{x}_1 + \mathbf{g})}{\mathbf{d}_1^T \mathbf{Q} \mathbf{d}_1} = -\frac{\begin{pmatrix} -8 & -3 \end{pmatrix} \begin{pmatrix} 8 \\ 3 \end{pmatrix}}{\begin{pmatrix} -8 & -3 \end{pmatrix} \begin{pmatrix} 4 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} -8 \\ -3 \end{pmatrix}} = \frac{73}{331}$$

- Compute the next iterate \mathbf{x}_2 :

$$\mathbf{x}_2 = \mathbf{x}_1 + \alpha_1 \mathbf{d}_1 = \begin{pmatrix} 2 \\ 1 \end{pmatrix} + \frac{73}{331} \begin{pmatrix} -8 \\ -3 \end{pmatrix} \approx \begin{pmatrix} 0.2356 \\ 0.3384 \end{pmatrix}$$

Example: Iteration 1 (Cont'd)

- For the stopping criterion, we need to compute $\|Q\mathbf{x}_2 + \mathbf{g}\|$:

$$Q\mathbf{x}_2 + \mathbf{g} = \begin{pmatrix} 4 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} 0.2356 \\ 0.3384 \end{pmatrix} + \begin{pmatrix} -1 \\ -2 \end{pmatrix} \approx \begin{pmatrix} 0.2810 \\ -0.7492 \end{pmatrix}$$

As $\|Q\mathbf{x}_2 + \mathbf{g}\| = 0.80 > \epsilon$, we continue the current iteration

- Compute β_2 :

- ▶ From the previous calculations: $Q\mathbf{x}_1 + \mathbf{g} = (8 \ 3)^T$

- ▶ Then:

$$\beta_2 = \frac{(Q\mathbf{x}_2 + \mathbf{g})^T (Q\mathbf{x}_2 + \mathbf{g})}{(Q\mathbf{x}_1 + \mathbf{g})^T (Q\mathbf{x}_1 + \mathbf{g})} = \frac{0.6401}{73} \approx 0.0088$$

- Compute the new direction \mathbf{d}_2 :

$$\mathbf{d}_2 = -Q\mathbf{x}_2 - \mathbf{g} + \beta_2 \mathbf{d}_1 = \begin{pmatrix} -0.2810 \\ 0.7492 \end{pmatrix} + \beta_2 \begin{pmatrix} -8 \\ -3 \end{pmatrix} = \begin{pmatrix} -0.3511 \\ -0.7229 \end{pmatrix}$$

Example: Iteration 2

- $k = 2$
- We compute α_2 :

$$\alpha_2 = -\frac{\mathbf{d}_2^T (\mathbf{Q}\mathbf{x}_2 + \mathbf{g})}{\mathbf{d}_2^T \mathbf{Q}\mathbf{d}_2} = 0.4122$$

- Next iterate \mathbf{x}_3 :

$$\mathbf{x}_3 = \mathbf{x}_2 + \alpha_2 \mathbf{d}_2 = \begin{pmatrix} 0.2356 \\ 0.3384 \end{pmatrix} + 0.4122 \begin{pmatrix} -0.3511 \\ -0.7229 \end{pmatrix} = \begin{pmatrix} 0.0909 \\ 0.6364 \end{pmatrix}$$

- Stopping criterion: $k = 2$ (= the dimension of the problem) and the process is **completed**

The CG algorithm has found the solution of this linear system in two iterations ! This point is the unique solution of our problem