

A short introduction to optimization

From unconstrained to constrained optimization

Gilles Gasso

INSA Rouen - ITI Departement
LITIS Laboratory

September 21, 2025

Plan

- 1 Unconstrained optimization
 - Formulation
 - Optimality conditions
 - Descent algorithms
 - Main methods
 - Determination of the step size
 - Illustration of descent methods
- 2 Constrained optimization
 - Formulation
 - Concept of Lagrangian and duality, condition of optimality
 - Lagrangian formulation
 - Optimality conditions
 - Duality and dual problem
 - Specific constrained optimization problems
- 3 Conclusion

Unconstrained optimization

Elements of the problem

- $\boldsymbol{\theta} \in \mathbb{R}^d$: vector of unknown real parameters
- $J : \mathbb{R}^d \rightarrow \mathbb{R}$: the function to be minimized
- Assumption: J is differentiable all over its domain
 $\text{dom} J = \{\boldsymbol{\theta} \in \mathbb{R}^d \mid J(\boldsymbol{\theta}) < \infty\}$

Problem formulation

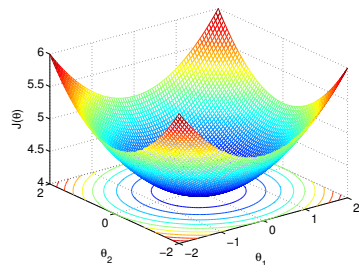
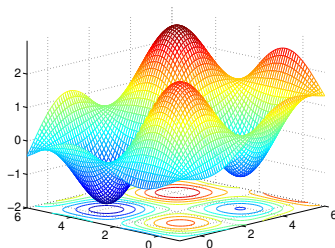
$$(P) \quad \min_{\boldsymbol{\theta} \in \mathbb{R}^d} J(\boldsymbol{\theta})$$

Unconstrained optimization

Examples

$$J(\theta) = \frac{1}{2} \theta^\top P \theta + q^\top \theta + r$$

with P a positive definite matrix



$$J(\theta) = \cos(\theta_1 - \theta_2) + \sin(\theta_1 + \theta_2) + \frac{\theta_1}{4}$$

Different solutions

Global solution

θ^* is said to be the global minimum solution of the problem if

$$J(\theta^*) \leq J(\theta), \quad \forall \theta \in \text{dom} J$$

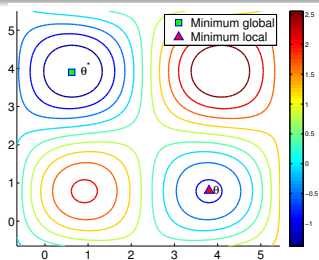
Local solution

$\hat{\theta}$ is a local minimum solution of problem (P) if it holds

$$J(\hat{\theta}) \leq J(\theta), \quad \forall \theta \in \text{dom} J \text{ such that } \|\hat{\theta} - \theta\| \leq \epsilon, \epsilon > 0$$

Illustration

$$J(\theta) = \cos(\theta_1 - \theta_2) + \sin(\theta_1 + \theta_2) + \frac{\theta_1}{4}$$



Optimality conditions

- How to assess a solution to the problem?

First order necessary condition

Theorem [First order condition]

Let $J : \mathbb{R}^d \rightarrow \mathbb{R}$ be a differential function on its domain. A vector $\hat{\theta}$ is a (local or global) solution of the problem (P), if it necessarily satisfies the condition $\nabla J(\hat{\theta}) = 0$.

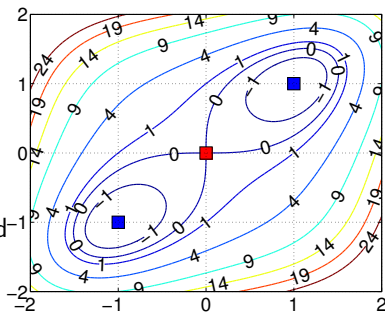
Remarks

- Any vector θ_0 that verifies $\nabla J(\theta_0) = 0$ is called a stationary point
- $\nabla J(\theta) \in \mathbb{R}^d$ is the gradient vector of J at θ .
- The gradient is the unique vector such that the directional derivative can be written as:

$$\lim_{t \rightarrow 0} \frac{J(\theta + t\mathbf{h}) - J(\theta)}{t} = \nabla J(\theta)^\top \mathbf{h}, \quad \mathbf{h} \in \mathbb{R}^d, \quad t \in \mathbb{R}$$

Example of a first order optimality condition

- $J(\boldsymbol{\theta}) = \theta_1^4 + \theta_2^4 - 4\theta_1\theta_2$
- Gradient $\nabla J(\boldsymbol{\theta}) = \begin{pmatrix} 4\theta_1^3 - 4\theta_2 \\ -4\theta_1 + 4\theta_2^3 \end{pmatrix}$
- Stationary points that verify $\nabla J(\boldsymbol{\theta}) = 0$.
- Three solutions $\boldsymbol{\theta}^{(1)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, $\boldsymbol{\theta}^{(2)} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ and $\boldsymbol{\theta}^{(3)} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$



Remarks

- $\boldsymbol{\theta}^{(2)}$ and $\boldsymbol{\theta}^{(3)}$ are local minimal but not $\boldsymbol{\theta}^{(1)}$
- every stationary point can be deemed a local extremum

We need another optimality condition

How to ensure that a stationary point is a minimum solution?

Hessian matrix

Twice differential function

$J : \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be a twice differentiable function on its domain $\text{dom}J$ if, at every point $\theta \in$, there exists a **unique symmetric matrix** $\mathbf{H}(\theta) \in \mathbb{R}^{d \times d}$ called **Hessian matrix** such that

$$J(\theta + \mathbf{h}) = J(\theta) + \nabla J(\theta)^\top \mathbf{h} + \mathbf{h}^\top \mathbf{H}(\theta) \mathbf{h} + \|\mathbf{h}\|^2 \varepsilon(\mathbf{h}).$$

$\varepsilon(\mathbf{h})$ is a continuous function at $\mathbf{0}$ with $\lim_{\mathbf{h} \rightarrow \mathbf{0}} \varepsilon(\mathbf{h}) = 0$

- $\mathbf{H}(\theta)$ is the second derivative matrix

$$\mathbf{H}(\theta) = \begin{pmatrix} \frac{\partial^2 J}{\partial \theta_1 \partial \theta_1} & \frac{\partial^2 J}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 J}{\partial \theta_1 \partial \theta_d} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{\partial^2 J}{\partial \theta_d \partial \theta_1} & \frac{\partial^2 J}{\partial \theta_d \partial \theta_2} & \cdots & \frac{\partial^2 J}{\partial \theta_d \partial \theta_d} \end{pmatrix}$$

- $\mathbf{H}(\theta) = \nabla_{\theta^\top} (\nabla_{\theta} J(\theta))$ is the Jacobian of the gradient function

Second order optimality condition

Theorem [Second order optimality condition]

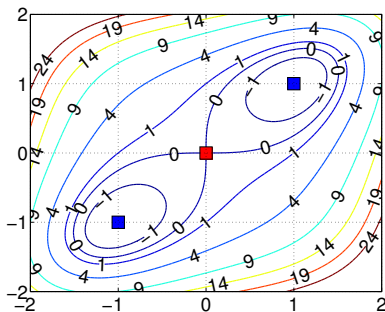
Let $J : \mathbb{R}^d \rightarrow \mathbb{R}$ be a twice differentiable function on its domain. If $\hat{\theta}$ is a minimum of J , then $\nabla J(\hat{\theta}) = 0$ and $\mathbf{H}(\hat{\theta})$ is a positive definite matrix.

Remarks

- \mathbf{H} is positive definite if and only if all its eigenvalues are positive
- \mathbf{H} is negative definite if and only if all its eigenvalues are negative
- For $\theta \in \mathbb{R}$, this condition means that the gradient of J at the minimum is null, $J'(\theta) = 0$ and its second derivative is positive i.e. $J''(\theta) > 0$
- If at a stationary point θ_0 , $\mathbf{H}(\hat{\theta})$ is negative definite, $\hat{\theta}$ is a local maximum of J

Illustration of the second order optimality condition

- $J(\boldsymbol{\theta}) = \theta_1^4 + \theta_2^4 - 4\theta_1\theta_2$
- Gradient : $\nabla J(\boldsymbol{\theta}) = \begin{pmatrix} 4\theta_1^3 - 4\theta_2 \\ -4\theta_1 + 4\theta_2^3 \end{pmatrix}$
- Stationary points : $\boldsymbol{\theta}^{(1)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, $\boldsymbol{\theta}^{(2)} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$
 $\boldsymbol{\theta}^{(3)} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$
- Hessian matrix $\mathbf{H}(\boldsymbol{\theta}) = \begin{pmatrix} 12\theta_1^2 & -4 \\ -4 & 12\theta_2^2 \end{pmatrix}$



	$\boldsymbol{\theta}^{(1)}$	$\boldsymbol{\theta}^{(2)}$	$\boldsymbol{\theta}^{(3)}$
Hessian	$\begin{pmatrix} 0 & -4 \\ -4 & 0 \end{pmatrix}$	$\begin{pmatrix} 12 & -4 \\ -4 & 12 \end{pmatrix}$	$\begin{pmatrix} 12 & -4 \\ -4 & 12 \end{pmatrix}$
Eigenvalues	4, -4	8, 16	8, 16
Type of solution	Saddle point	Minimum	Minimum

Necessary and sufficient optimality condition

Theorem [2nd order sufficient condition]

Assume the hessian matrix $\mathbf{H}(\hat{\boldsymbol{\theta}})$ of $J(\boldsymbol{\theta})$ at $\hat{\boldsymbol{\theta}}$ exists and is positive definite. Assume also the gradient $\nabla J(\hat{\boldsymbol{\theta}}) = 0$. Then $\hat{\boldsymbol{\theta}}$ is a (local or global) minimum of problem (P).

Theorem [Sufficient and necessary optimality condition]

Let J be a convex function. Every local solution $\hat{\boldsymbol{\theta}}$ is a global solution $\boldsymbol{\theta}^$.*

Recall

A function $J : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if it verifies

$$J(\alpha \boldsymbol{\theta} + (1 - \alpha) \mathbf{z}) \leq \alpha J(\boldsymbol{\theta}) + (1 - \alpha) J(\mathbf{z}), \quad \forall \boldsymbol{\theta}, \mathbf{z} \in \text{dom} J, \quad 0 \leq \alpha \leq 1$$

How to find the solution(s)?

- We have seen how to assess a solution to the problem
- Now, how to compute a solution?

Principle of descent algorithms

Direction of descent

Let the function $J : \mathbb{R}^d \rightarrow \mathbb{R}$. The vector $\mathbf{h} \in \mathbb{R}^d$ is called a **descent direction** in $\boldsymbol{\theta}$ if there exists $\alpha > 0$ such that $J(\boldsymbol{\theta} + \alpha \mathbf{h}) < J(\boldsymbol{\theta})$

Principle of descent methods

- Start from an initial point $\boldsymbol{\theta}_0$
 - Design a sequence of points $\{\boldsymbol{\theta}_k\}$ with $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha_k \mathbf{h}_k$
 - Ensure that the sequence $\{\boldsymbol{\theta}_k\}$ converges to a stationary point $\hat{\boldsymbol{\theta}}$
-
- \mathbf{h}_k : direction of descent
 - α_k : **step size**

General approach

General algorithm

- 1: Let $k = 0$, initialize θ_k
- 2: **repeat**
- 3: Find a descent direction $\mathbf{h}_k \in \mathbb{R}^d$
- 4: Line search: find a step size $\alpha_k > 0$ in the direction \mathbf{h}_k such that $J(\theta_k + \alpha_k \mathbf{h}_k)$ decreases "enough"
- 5: Update: $\theta_{k+1} \leftarrow \theta_k + \alpha_k \mathbf{h}_k$ and $k \leftarrow k + 1$
- 6: **until** convergence

- The methods of descent differ by the choice of:
 - \mathbf{h} : gradient algorithm, Newton, Quasi-Newton algorithm
 - α : backtracking...

Gradient Algorithm

Theorem [descent direction and opposite direction of gradient]

Let $J(\boldsymbol{\theta})$ be a differential function. The direction $\mathbf{h} = -\nabla J(\boldsymbol{\theta}) \in \mathbb{R}^d$ is a descent direction.

Proof.

J being differentiable, for any $t > 0$ we have

$J(\boldsymbol{\theta} + t\mathbf{h}) = J(\boldsymbol{\theta}) + t\nabla J(\boldsymbol{\theta})^\top \mathbf{h} + t\|\mathbf{h}\|\epsilon(t\mathbf{h})$. Setting $\mathbf{h} = -\nabla J(\boldsymbol{\theta})$, we get $J(\boldsymbol{\theta} + t\mathbf{h}) - J(\boldsymbol{\theta}) = -t\|\nabla J(\boldsymbol{\theta})\|^2 + t\|\mathbf{h}\|\epsilon(t\mathbf{h})$. For t small enough $\epsilon(t\mathbf{h}) \rightarrow 0$ and so $J(\boldsymbol{\theta} + t\mathbf{h}) - J(\boldsymbol{\theta}) = -t\|\nabla J(\boldsymbol{\theta})\|^2 < 0$. It is then a descent direction. \square

Characteristics of the gradient algorithm

- Choice of the descent direction at $\boldsymbol{\theta}_k$: $\mathbf{h}_k = -\nabla J(\boldsymbol{\theta}_k)$
- Complexity of the update: $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k - \alpha_k \nabla J(\boldsymbol{\theta}_k)$ costs $\mathcal{O}(d)$

Newton algorithm

- 2nd order approximation of J at $\boldsymbol{\theta}_k$

$$J(\boldsymbol{\theta} + \mathbf{h}) \approx J(\boldsymbol{\theta}_k) + \nabla J(\boldsymbol{\theta}_k)^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{H}(\boldsymbol{\theta}_k) \mathbf{h}$$

with $\mathbf{H}(\boldsymbol{\theta}_k)$ the positive definite Hessian matrix

- The direction \mathbf{h}_k which minimizes this approximation is obtained by

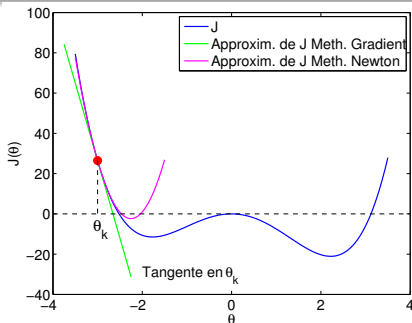
$$\nabla J(\boldsymbol{\theta} + \mathbf{h}_k) = 0 \quad \Rightarrow \quad \mathbf{h}_k = -\mathbf{H}(\boldsymbol{\theta}_k)^{-1} \nabla J(\boldsymbol{\theta}_k)$$

Features

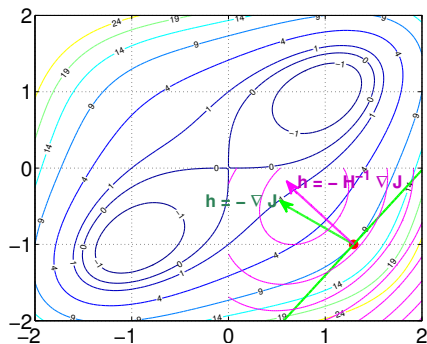
- Descent direction at $\boldsymbol{\theta}_k$: $\mathbf{h}_k = -\mathbf{H}(\boldsymbol{\theta}_k)^{-1} \nabla J(\boldsymbol{\theta}_k)$
- Complexity of the update: $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k - \alpha_k \mathbf{H}(\boldsymbol{\theta}_k)^{-1} \nabla J(\boldsymbol{\theta}_k)$ costs $\mathcal{O}(d^3)$ flops
- $\mathbf{H}(\boldsymbol{\theta}_k)$ is not always guaranteed to be positive definite matrix. Hence we cannot always ensure that \mathbf{h}_k is a direction of descent

Illustration of gradient and Newton methods

Local approximation of the two methods in 1D



Directions of descent in 2D



Set up the step size α_k in the update $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha_k \mathbf{h}_k$

- Fixed step size: use a fixed value $\alpha_k = \alpha > 0$ at each iteration k
- Variable step size: α_k is adaptative using a line search

Armijo's rule: choose α_k in order to have a sufficient decrease of J i.e.

$$J(\boldsymbol{\theta}_k + \alpha_k \mathbf{h}) \leq J(\boldsymbol{\theta}_k) + c \alpha_k \nabla J(\boldsymbol{\theta}_k)^\top \mathbf{h}_k$$

- Usually c is chosen in the range $[10^{-5}, 10^{-1}]$
- \mathbf{h}_k is a descent direction, we have $\nabla J(\boldsymbol{\theta}_k)^\top \mathbf{h}_k < 0$, thus the decrease of J

Backtracking

- 1: Fix an initial step $\bar{\alpha}$, choose $0 < \rho < 1$, $\alpha \leftarrow \bar{\alpha}$
- 2: **repeat**
- 3: $\alpha \leftarrow \rho \alpha$
- 4: **until** $J(\boldsymbol{\theta}_k + \alpha \mathbf{h}) > J(\boldsymbol{\theta}_k) + c \alpha \nabla J(\boldsymbol{\theta}_k)^\top \mathbf{h}_k$

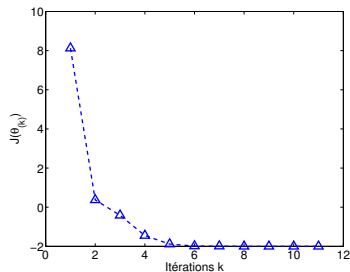
Choice of the initial step

- Newton method:
 $\bar{\alpha} = 1$
- Gradient method:
 $\bar{\alpha} = 2 \frac{J(\boldsymbol{\theta}_k) - J(\boldsymbol{\theta}_{k-1})}{\nabla J(\boldsymbol{\theta}_k)^\top \mathbf{h}_k}$

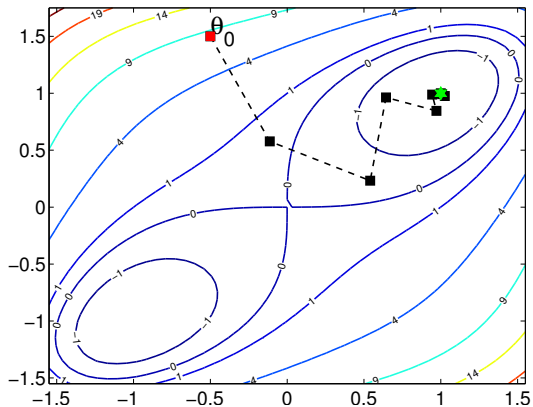
Interpretation: as long as J does not decrease, the step size is decreased

Gradient method

J along the iterations

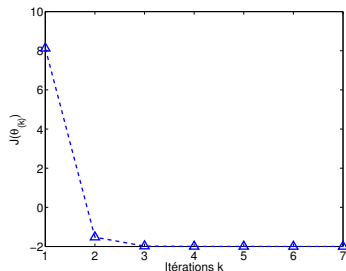


Evolution of the iterates

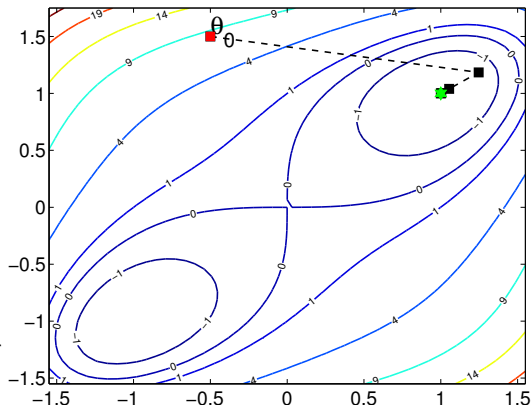


Newton method

J along the iterations



Evolution of the iterates

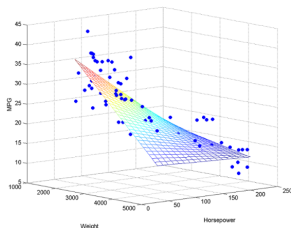


- At each iteration we considered the matrix $\mathbf{H}(\theta) + \lambda \mathbf{I}$ instead of \mathbf{H} to guarantee the positive definite property of Hessian

Constrained optimization problems

Examples and formulation

Example 1: sparse Regression



- Output to be predicted: $y \in \mathbb{R}$
- Input variables: $\mathbf{x} \in \mathbb{R}^d$
- Linear model: $f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\theta}$
- $\boldsymbol{\theta} \in \mathbb{R}^d$: parameters of the model

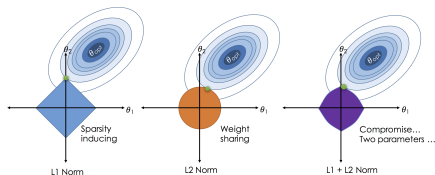
Determination of a sparse $\boldsymbol{\theta}$

- Minimization of square error
- Only a few parameters are non-zero

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \boldsymbol{\theta})^2$$

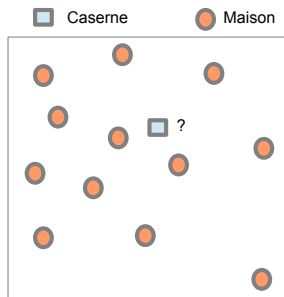
s.t. $\|\boldsymbol{\theta}\|_p \leq k$

with $\|\boldsymbol{\theta}\|_p^p = \sum_{j=1}^d |\theta_j|^p$



http://www.ds100.org/sp17/assets/notebooks/linear_regression/Regularization.html

Example 2: where to settle the firehouse?



- House M_i : defined by its coordinates $\mathbf{z}_i = [x_i, y_i]^\top$
- Let $\boldsymbol{\theta}$ be the coordinates of the firehouse
- Minimize the distance from the firehouse to the farthest house

Problem formulation

$$\min_{\boldsymbol{\theta}} \max_{i=1, \dots, n} \|\boldsymbol{\theta} - \mathbf{z}_i\|^2$$

Equivalent problem

$$\begin{aligned} & \min_{t \in \mathbb{R}, \boldsymbol{\theta} \in \mathbb{R}^2} t \\ & \text{s.t. } \|\boldsymbol{\theta} - \mathbf{z}_i\|^2 \leq t \quad \forall i = 1, \dots, n \end{aligned}$$

Formulation of constrained optimization problem

Notations and assumptions

- $\theta \in \mathbb{R}^d$: vector of unknown real parameters
- $J : \mathbb{R}^d \rightarrow \mathbb{R}$, the function to be minimized on its domain $\text{dom} J$
- f_i and g_j are differentiable functions of \mathbb{R}^d on \mathbb{R}

Primal problem \mathcal{P}

$$\begin{array}{ll}
 \min_{\theta \in \mathbb{R}^d} & J(\theta) \\
 \text{s.t.} & f_i(\theta) = 0 \quad \forall i = 1, \dots, n \\
 & g_j(\theta) \leq 0 \quad \forall j = 1, \dots, m
 \end{array}$$

objective function
 n Equality Constraints
 m Inequality Constraints

Feasibility

Let $p^* = \min_{\theta} \{J(\theta) \text{ such that } f_i(\theta) = 0 \forall i \text{ and } g_j(\theta) \leq 0 \forall j\}$

- If $p^* = \infty$ then the problem does not admit a feasible solution

Characterization of the solutions

Feasibility domain

The feasible domain is defined by the set of constraints

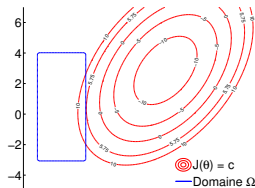
$$\Omega(\boldsymbol{\theta}) = \left\{ \boldsymbol{\theta} \in \mathbb{R}^d; \ f_i(\boldsymbol{\theta}) = 0 \ \forall i \text{ and } g_j(\boldsymbol{\theta}) \leq 0 \ \forall j \right\}$$

Feasible points

- *$\boldsymbol{\theta}_0$ is feasible if $\boldsymbol{\theta}_0 \in \text{dom}J$ and $\boldsymbol{\theta}_0 \in \Omega(\boldsymbol{\theta})$ ie $\boldsymbol{\theta}_0$ fulfills all the constraints and $J(\boldsymbol{\theta}_0)$ has a finite value*
- *$\boldsymbol{\theta}^*$ is a global solution of the problem if $\boldsymbol{\theta}^*$ is a feasible solution such that $J(\boldsymbol{\theta}^*) \leq J(\boldsymbol{\theta})$ for every $\boldsymbol{\theta}$*
- *$\hat{\boldsymbol{\theta}}$ is a local optimal solution if $\hat{\boldsymbol{\theta}}$ is feasible and $J(\hat{\boldsymbol{\theta}}) \leq J(\boldsymbol{\theta})$ for every $\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\| \leq \epsilon$*

Example 1

$$\begin{aligned}
 \min_{\boldsymbol{\theta}} \quad & 0.9\theta_1^2 - 0.74\theta_1\theta_2 \\
 & + 0.75\theta_2^2 - 5.4\theta_1 - 1.2\theta_2 \\
 \text{s.t.} \quad & -4 \leq \theta_1 \leq -1 \\
 & -3 < \theta_2 < 4
 \end{aligned}$$



- Parameters: $\boldsymbol{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}$
- Objective function:

$$J(\boldsymbol{\theta}) = 0.9\theta_1^2 - 0.74\theta_1\theta_2 + 0.75\theta_2^2 - 5.4\theta_1 - 1.2\theta_2$$

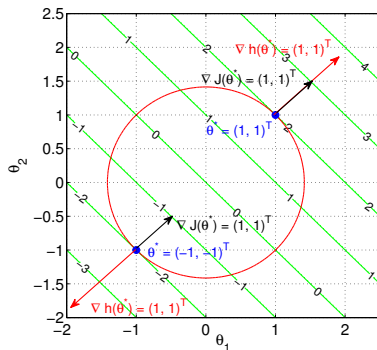
- Feasibility domain (four inequality constraints):

$$\Omega(\boldsymbol{\theta}) = \{\boldsymbol{\theta} \in \mathbb{R}^2; -4 \leq \theta_1 \leq -1 \text{ and } -3 \leq \theta_2 \leq 4\}$$

Example 2

Example

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & \theta_1 + \theta_2 \\ \text{s.t.} \quad & \theta_1^2 + \theta_2^2 - 2 = 0 \end{aligned}$$



- An equality constraint
- Domain of feasibility: a circle with center at $\mathbf{0}$ and diameter equals to 2
- The optimal solution is obtained for $\theta^* = (-1 \quad -1)^\top$ and we have $J(\theta^*) = -2$

Optimality

- How to assess a solution of the primal problem?
- Do we have optimality conditions similar to those of unconstrained optimization?

Notion of Lagrangian

Primal problem \mathcal{P}

$$\begin{array}{ll}
 \min_{\boldsymbol{\theta} \in \mathbb{R}^d} & J(\boldsymbol{\theta}) \\
 \text{s.t.} & f_i(\boldsymbol{\theta}) = 0 \quad \forall i = 1, \dots, n \\
 & g_j(\boldsymbol{\theta}) \leq 0 \quad \forall j = 1, \dots, m
 \end{array}$$

n equality constraints
 m inequality constraints

$\boldsymbol{\theta}$ is called primal variable

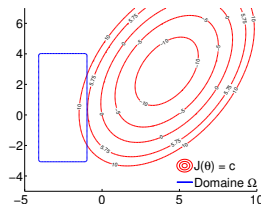
Principle of Lagrangian

- Each constraint is associated to a scalar parameter called **Lagrange multiplier**
- Equality constraint** $f_i(\boldsymbol{\theta}) = 0$: we associate $\lambda_i \in \mathbb{R}$
- Inequality constraint** $g_j(\boldsymbol{\theta}) \leq 0$: we associate $\alpha_j \geq 0^a$
- Lagrangian allows to transform the problem with constraints into a problem without constraints with additional variables: λ_i and α_j .

^aBeware of the type of inequality i.e. $g_j(\boldsymbol{\theta}) \leq 0$

Example

$$\begin{aligned}
 \min_{\theta \in \mathbb{R}^2} \quad & 0.9\theta_1^2 - 0.74\theta_1\theta_2 \\
 & + 0.75\theta_1^2 - 5.4\theta_1 - 1.2\theta_2 \\
 \text{s.t.} \quad & -4 \leq \theta_1 \leq -1 \\
 & -3 \leq \theta_2 \leq 4
 \end{aligned}$$



Constraints (inequality)

- ① $-4 \leq \theta_1 \Leftrightarrow -\theta_1 - 4 \leq 0$
- ② $\theta_1 \leq -1 \Leftrightarrow \theta_1 + 1 \leq 0$
- ③ $-3 \leq \theta_2 \Leftrightarrow -\theta_2 - 3 \leq 0$
- ④ $\theta_2 \leq 4 \Leftrightarrow -\theta_2 - 4 \leq 0$

Related Lagrange Parameters

- ① $\alpha_1 \geq 0$
- ② $\alpha_2 \geq 0$
- ③ $\alpha_3 \geq 0$
- ④ $\alpha_4 \geq 0$

Lagrangian

$$\begin{array}{ll} \min_{\boldsymbol{\theta} \in \mathbb{R}^d} & J(\boldsymbol{\theta}) \\ & f_i(\boldsymbol{\theta}) = 0 \quad \forall i = 1, \dots, n \\ \text{s.c.} & g_j(\boldsymbol{\theta}) \leq 0 \quad \forall j = 1, \dots, m \end{array}$$

Associated Lagrange parameters

None

λ_i any real number $\forall i = 1, \dots, n$

$\alpha_j \geq 0 \quad \forall j = 1, \dots, m$

Lagrangian

The Lagrangian is defined by :

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = J(\boldsymbol{\theta}) + \sum_{i=1}^n \lambda_i f_i(\boldsymbol{\theta}) + \sum_{j=1}^m \alpha_j g_j(\boldsymbol{\theta}) \quad \text{avec} \quad \mu_j \geq 0, \forall j = 1, \dots, m$$

- Lagrange parameters $\lambda_i, i = 1, \dots, n$ and $\alpha_j, j = 1, \dots, m$ are called **dual variables**
- **Dual variables are unknown parameters** to be determined

Examples

Example 1

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & 0.9\theta_1^2 - 0.74\theta_1\theta_2 + 0.75\theta_1^2 - 5.4\theta_1 - 1.2\theta_2 \\ \text{s.t.} \quad & -4 \leq \theta_1 \leq -1 \quad \text{and} \quad -3 \leq \theta_2 \leq 4 \end{aligned}$$

Lagrangian

$$\begin{aligned} \mathcal{L}(\alpha, \theta) = & 0.9\theta_1^2 - 0.74\theta_1\theta_2 + 0.75\theta_1^2 - 5.4\theta_1 - 1.2\theta_2 \\ & + \alpha_1(-\theta_1 - 4) + \alpha_2(\theta_1 + 1) + \alpha_3(-\theta_2 - 3) + \alpha_4(-\theta_2 - 4) \end{aligned}$$

with $\alpha_1 \geq 0, \alpha_2 \geq 0, \alpha_3 \geq 0, \alpha_4 \geq 0$ (because of inequality constraints)

Example 2

$$\begin{aligned} \min_{\theta \in \mathbb{R}^3} \quad & \frac{1}{2} (\theta_1^2 + \theta_2^2 + \theta_3^2) \\ \text{s.t.} \quad & \theta_1 + \theta_2 + 2\theta_3 = 1 \quad \text{equality constraint} \\ & \theta_1 + 4\theta_2 + 2\theta_3 = 3 \quad \text{equality constraint} \end{aligned}$$

Lagrangian

$$\mathcal{L}(\lambda, \theta) = \frac{1}{2} (\theta_1^2 + \theta_2^2 + \theta_3^2) + \lambda_1(\theta_1 + \theta_2 + 2\theta_3 - 1) + \lambda_2(\theta_1 + 4\theta_2 + 2\theta_3 - 3)$$

with $\lambda_1, \lambda_2 \in \mathbb{R}$ (equality constraints)

Necessary optimality conditions

Assume that J, f_i, g_j are differentiable functions. Let θ^* be a feasible solution to the problem \mathcal{P} . Then there exists dual variables $\lambda_i^*, i = 1, \dots, n, \alpha_j^*, j = 1, \dots, m$ such that the KKT conditions are met.

Karush-Kuhn-Tucker (KKT) Conditions

Stationarity

$$\begin{aligned} \nabla \mathcal{L}(\lambda, \alpha, \theta) &= 0 \quad \text{ie} \\ \nabla J(\theta) + \sum_{i=1}^n \lambda_i \nabla f_i(\theta) + \sum_{j=1}^m \alpha_j \nabla g_j(\theta) &= 0 \end{aligned}$$

Primal feasibility

$$\begin{aligned} f_i(\theta) &= 0 & \forall i = 1, \dots, n \\ g_j(\theta) &\leq 0 & \forall j = 1, \dots, m \end{aligned}$$

Dual feasibility

$$\alpha_j \geq 0 \quad \forall j = 1, \dots, m$$

Complementary slackness

$$\alpha_j g_j(\theta) = 0 \quad \forall j = 1, \dots, m$$

Example

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & \frac{1}{2}(\theta_1^2 + \theta_2^2) \\ \text{s.t.} \quad & \theta_1 - 2\theta_2 + 2 \leq 0 \end{aligned}$$

- Lagrangian : $\mathcal{L}(\alpha, \theta) = \frac{1}{2}(\theta_1^2 + \theta_2^2) + \alpha(\theta_1 - 2\theta_2 + 2)$, $\alpha \geq 0$

- KKT Conditions

- Stationarity: $\nabla_{\theta} \mathcal{L}(\alpha, \theta) = 0 \quad \Rightarrow \quad \begin{cases} \theta_1 = -\alpha \\ \theta_2 = -2\alpha \end{cases}$

- Primal feasibility : $\theta_1 - 2\theta_2 + 2 \leq 0$

- Dual feasibility : $\alpha \geq 0$

- Complementary slackness : $\alpha(\theta_1 - 2\theta_2 + 2) = 0$

- Remarks on the complementary slackness

- If $\theta_1 - 2\theta_2 + 2 < 0$ (inactive constraint) $\Rightarrow \alpha = 0$ (no penalty required as the constraint is satisfied)
- If $\mu > 0 \Rightarrow \theta_1 - 2\theta_2 + 2 = 0$ (active constraint)

Duality

Dual function

Let $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\lambda}, \boldsymbol{\alpha})$ be the lagrangian of the primal problem \mathcal{P} with $\alpha_j \geq 0$.
The corresponding *dual function* is defined as

$$\mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\alpha}) = \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

Theorem [Weak duality]

Let $p^* = \min_{\boldsymbol{\theta}} \{J(\boldsymbol{\theta}) \text{ such that } f_i(\boldsymbol{\theta}) = 0 \forall i \text{ and } g_j(\boldsymbol{\theta}) \leq 0 \forall j\}$ be the optimum value (supposed finite) of the problem \mathcal{P} . Then, for any value of $\alpha_j \geq 0, \forall j$ and $\lambda_i, \forall i$, we have

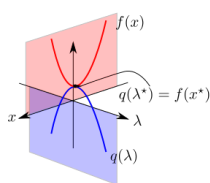
$$\mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu}) \leq p^*$$

Dual problem

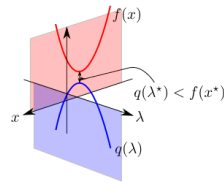
- The weak duality indicates that the dual function $\mathcal{D}(\lambda, \alpha) = \min_{\theta} \mathcal{L}(\theta, \lambda, \alpha)$ is a lower bound of p^*
- Bridge the gap: maximize the dual w.r.t. dual variables λ and μ to make this lower bound close to p^*

Dual problem

$$\begin{aligned} \max_{\lambda, \alpha} \quad & \mathcal{D}(\lambda, \mu) \\ \text{s.t.} \quad & \alpha_j \geq 0 \quad \forall j = 1, \dots, m \end{aligned}$$



strong duality



weak duality

<http://www.onmyphd.com/?p=duality.theory>

Interest of the dual problem

Remarks

- Transform the primal problem into an equivalent dual problem possibly much simpler to solve
- Solving the dual problem can lead to the solution of the primal problem
- Solving the dual problem gives the optimal values of the Lagrange multipliers

Example : inequality constraints

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & \frac{1}{2}(\theta_1^2 + \theta_2^2) \\ \text{s.t.} \quad & \theta_1 - 2\theta_2 + 2 \leq 0 \end{aligned}$$

- Lagrangian : $\mathcal{L}(\theta, \alpha) = \frac{1}{2}(\theta_1^2 + \theta_2^2) + \alpha(\theta_1 - 2\theta_2 + 2)$, $\alpha \geq 0$

- Stationarity of the KKT Condition :

$$\nabla_{\theta} \mathcal{L}(\mu, \theta) = 0 \quad \Rightarrow \quad \begin{cases} \theta_1 = -\alpha \\ \theta_2 = 2\alpha \end{cases} \quad (1)$$

- Dual function $\mathcal{D}(\alpha) = \min_{\theta} L(\theta, \alpha)$: by substituting (1) in \mathcal{L} we obtain

$$\mathcal{D}(\alpha) = -\frac{5}{2}\alpha^2 + 2\alpha$$

- Dual problem : $\max_{\alpha} \mathcal{D}(\alpha)$ s.c. $\alpha \geq 0$

- Dual solution

$$\nabla \mathcal{D}(\alpha) = 0 \Rightarrow \alpha = \frac{2}{5} \quad (\text{that satisfies } \alpha \geq 0) \quad (2)$$

- Primal solution : (2) and (1) lead to $\theta = \left(-\frac{2}{5} \quad \frac{4}{5}\right)^{\top}$

Convex constrained optimization

$$\begin{array}{ll} \min_{\boldsymbol{\theta} \in \mathbb{R}^d} & J(\boldsymbol{\theta}) \\ \text{s.t.} & f_i(\boldsymbol{\theta}) = 0 \quad \forall i = 1, \dots, n \\ & g_j(\boldsymbol{\theta}) \leq 0 \quad \forall j = 1, \dots, m \end{array}$$

Convexity condition

J is a convex function

f_i are linear $\forall i = 1, n$

g_j are convex functions $\forall j = 1, m$

Problems of interest

- Linear Programming (LP)
- Quadratic Programming (QP)
- Off-the-shelves toolboxes exist for those problems (Gurobi, Mosek, CVX ...)



QP convex problem

Standard form

$$\begin{array}{ll}
 \min_{\boldsymbol{\theta} \in \mathbb{R}^d} & \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{G} \boldsymbol{\theta} + \mathbf{q}^\top \boldsymbol{\theta} + r \\
 \text{s.t.} & \mathbf{a}_i^\top \boldsymbol{\theta} = b_i \quad \forall i = 1, \dots, n \quad \text{affine equality constraint} \\
 & \mathbf{c}_j^\top \boldsymbol{\theta} \geq d_j \quad \forall j = 1, \dots, m \quad \text{linear inequality constraints}
 \end{array}$$

with $\mathbf{q}, \mathbf{a}_i, \mathbf{c}_j \in \mathbb{R}^d$, b_i and d_j real scalar values and $\mathbf{G} \in \mathbb{R}^{d \times d}$ a **positive definite matrix**

Examples

SVM Problem

$$\begin{array}{ll}
 \min_{\boldsymbol{\theta} \in \mathbb{R}^2} & \frac{1}{2} (\theta_1^2 + \theta_2^2) \\
 \text{s.t.} & \theta_1 - 2\theta_2 + 2 \leq 0
 \end{array}
 \qquad
 \begin{array}{ll}
 \min_{\boldsymbol{\theta}, b} & \frac{1}{2} \|\boldsymbol{\theta}\|^2 \\
 \text{s.t.} & y_i (\boldsymbol{\theta}^\top \mathbf{x}_i + b) \geq 1 \quad \forall i = 1, N
 \end{array}$$

Conclusion

- Unconstrained optimization of smooth objective function
 - Characterization of the solution(s) requires checking the optimality conditions
 - Computation of a solution using descent methods
 - Gradient descent method
 - Newton method
 - Optimization under constraints
 - Lagrangian: allows to reduce to an unconstrained problem via Lagrange multipliers
 - To each constraint corresponds a multiplier \Rightarrow Lagrange parameters act as a penalty if the corresponding constraints are violated
 - Optimally (KKT conditions): Stationary condition + feasibility conditions + Complementary conditions
 - Duality: provides lower bound on the primal problem. Dual problem sometimes easier to solve than primal.