

# Descent methods for unconstrained optimization

Gilles Gasso

INSA Rouen - ASI Departement  
Laboratory LITIS

September 19, 2019

# Plan

- 1 Formulation
- 2 Optimality conditions
- 3 Descent algorithms
  - Main methods of descent
  - Research of the step
  - Summary
- 4 Illustration of descent methods

# Unconstrained optimization

## Elements of the problem

- $\theta \in \mathbb{R}^d$  : 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  
 $\text{dom}J = \{\theta \in \mathbb{R}^d \mid J(\theta) < \infty\}$

## Problem formulation

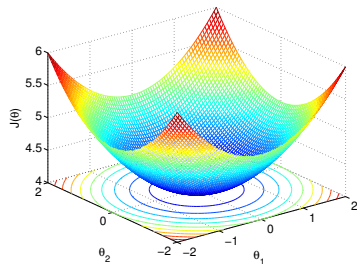
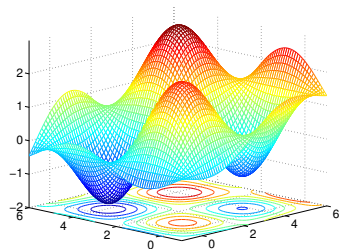
$$(P) \quad \min_{\theta \in \mathbb{R}^d} J(\theta)$$

# Unconstrained optimization

## Examples

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

with  $\mathbf{P}$  a positive definite matrix



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

# Different solutions

## Global solution

$\theta^*$  is said to be the global minimum solution of the problem if

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

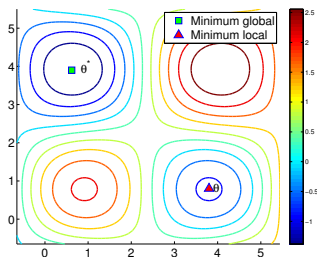
## Local solution

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

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

## Illustration

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



# Optimality conditions

- How do we 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  $\theta_0$  is a (local or global) solution of the problem (P), if it necessarily satisfies the condition  $\nabla J(\theta_0) = 0$ .

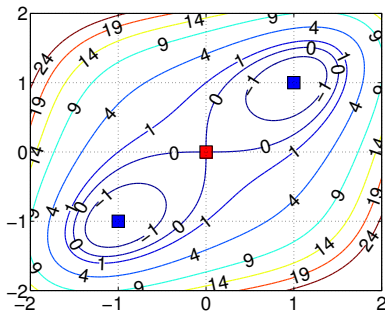
## Vocabulary

- Any vector  $\theta_0$  that verifies  $\nabla J(\theta_0) = 0$  is called a stationary point or critical point
- $\nabla J(\theta) \in \mathbb{R}^d$  is the gradient vector of  $J$  at  $\theta$ .
- 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

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



### Remarks

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

### We need another optimality condition

How to ensure that a stationary point is a minimum?



# 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  $\theta \in \text{dom}J$ , there exists **an unique symmetrical matrix  $H(\theta) \in \mathbb{R}^{d \times d}$**  called **Hessian matrix** such that

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

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

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

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

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

## Examples

## Example 1

- Objective function

$$J(\boldsymbol{\theta}) = \theta_1^4 + \theta_2^4 - 4\theta_1\theta_2$$

- Gradient

$$\nabla J(\boldsymbol{\theta}) = \begin{pmatrix} 4\theta_1^3 - 4\theta_2 \\ -4\theta_1 + 4\theta_2^3 \end{pmatrix}$$

- Hessian matrix

$$\mathbf{H}(\boldsymbol{\theta}) = \begin{pmatrix} 12\theta_1^2 & -4 \\ -4 & 12\theta_2^2 \end{pmatrix}$$

## Example 2

- Quadratic objective function

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

- Directional derivative

$$D(\mathbf{h}, \boldsymbol{\theta}) = \lim_{t \rightarrow 0} \frac{J(\boldsymbol{\theta} + t\mathbf{h}) - J(\boldsymbol{\theta})}{t}$$

$$D(\mathbf{h}, \boldsymbol{\theta}) = (\mathbf{P}\boldsymbol{\theta} + \mathbf{q})^\top \mathbf{h}$$

- Gradient  $\nabla J(\boldsymbol{\theta}) = \mathbf{P}\boldsymbol{\theta} + \mathbf{q}$

- Hessian matrix  $\mathbf{H}(\boldsymbol{\theta}) = \mathbf{P}$

## Second order optimality condition

### Theorem [Second order optimality condition]

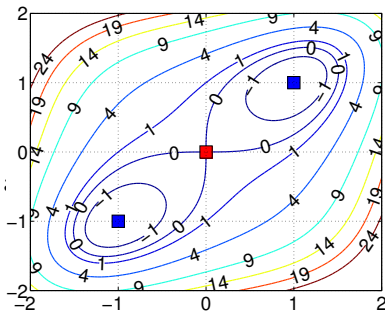
Let  $J : \mathbb{R}^d \rightarrow \mathbb{R}$  be a twice differentiable function on its domain. If  $\theta_0$  is a minimum of  $J$ , then  $\nabla J(\theta_0) = 0$  and  $\mathbf{H}(\theta_0)$  is a positive definite matrix.

### 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  $\theta_0$   $\mathbf{H}(\theta_0)$  is negative definite,  $\theta_0$  is a local maximum of  $J$

## Illustration of the second order optimality condition

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



	$\theta^{(1)}$	$\theta^{(2)}$	$\theta^{(3)}$
Hessian	$\begin{pmatrix} 0 & -4 \\ -4 & 12 \end{pmatrix}$	$\begin{pmatrix} 12 & -4 \\ -4 & 12 \end{pmatrix}$	$\begin{pmatrix} 12 & -4 \\ -4 & 12 \end{pmatrix}$
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  $\mathbf{H}(\boldsymbol{\theta}_0)$  of  $J(\boldsymbol{\theta})$  at  $\boldsymbol{\theta}_0$  exists and is positive definite. Assume also the gradient  $\nabla J(\boldsymbol{\theta}_0) = 0$ . Then  $\boldsymbol{\theta}_0$  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
- A question to be addressed is: 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 *direction of descent* 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
  - $\alpha_k$  is the **step size**

# General approach

## General algorithm

- 1: Let  $k = 0$ , initialize  $\boldsymbol{\theta}_k$
- 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  $J(\boldsymbol{\theta}_k + \alpha_k \mathbf{h}_k)$  decreases "enough"
- 5: Update:  $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha_k \mathbf{h}_k$  and  $k \leftarrow k + 1$
- 6: **until**  $\|\nabla J(\boldsymbol{\theta}_k)\| < \epsilon$

- The methods of descent differ by the choice of:
  - $\mathbf{h}$ : gradient algorithm, Newton, Quasi-Newton algorithm
  - $\alpha$ : Line search, 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 the twice differentiable  $J$  at  $\theta_k$

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

with  $\mathbf{H}(\theta_k)$  the positive definite Hessian matrix

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

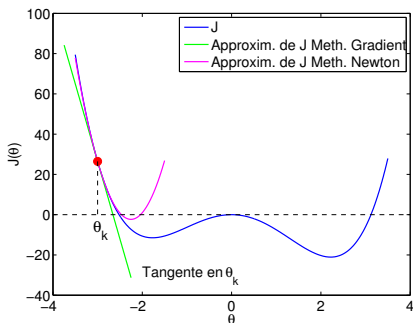
$$\nabla J(\theta + \mathbf{h}_k) = 0 \quad \Rightarrow \quad \mathbf{h}_k = -\mathbf{H}(\theta_k)^{-1} \nabla J(\theta_k)$$

### Features

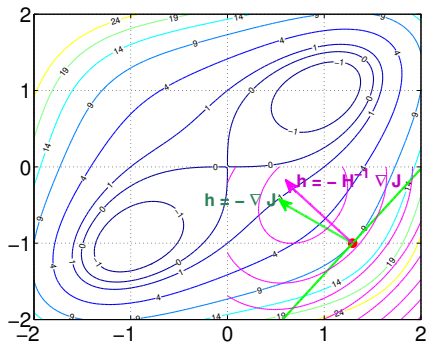
- Choice of the **descent direction at  $\theta_k$** :  $\mathbf{h}_k = -\mathbf{H}(\theta_k)^{-1} \nabla J(\theta_k)$
- **Complexity of the update**:  $\theta_{k+1} \leftarrow \theta_k - \alpha_k \mathbf{H}(\theta_k)^{-1} \nabla J(\theta_k)$  costs  $\mathcal{O}(d^3)$  flops
- $\mathbf{H}(\theta_k)$  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



# Quasi-Newton method

## Main features

- Choice of the **descent direction** at  $\theta_k$ :  $\mathbf{h}_k = -\mathbf{B}(\theta_k)^{-1}\nabla J(\theta_k)$
- $\mathbf{B}(\theta_k)$  is an positive definite approximation of the Hessian matrix
- **Complexity of the update**: most of the times  $\mathcal{O}(d^2)$

## Line search

Assume the direction of descent  $\mathbf{h}_k$  at  $\boldsymbol{\theta}_k$  is fixed. We aim to find the step size  $\alpha_k > 0$  in the direction  $\mathbf{h}_k$  such that the function  $J(\boldsymbol{\theta}_k + \alpha_k \mathbf{h}_k)$  decreases enough (compared to  $J(\boldsymbol{\theta}_k)$ )

### Several options

- Fixed step size: use a fixed value  $\alpha > 0$  at each iteration  $k$

$$\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha \mathbf{h}_k$$

- Optimal step size  $\alpha_k^*$

$$\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha_k^* \mathbf{h}_k \quad \text{with} \quad \alpha_k^* = \arg \min_{\alpha > 0} J(\boldsymbol{\theta}_k + \alpha \mathbf{h}_k)$$

- Variable step size: the choice  $\alpha_k$  is adapted to the current iteration

$$\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha_k \mathbf{h}_k$$

# Line search

## Pros and cons

- Fixed step size strategy: often not very effective
- Optimal step size: can be costly in calculation time
- Variable step: most commonly used approach
  - The step is often imprecise
  - A trade-off between computation cost and decrease of  $J$

## Variable step size

### Armijo's rule

Determine the step size  $\alpha_k$  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  $[10^{-5}, 10^{-1}]$
- Having  $\mathbf{h}_k$  the direction of descent, we have  $\nabla J(\boldsymbol{\theta}_k)^\top \mathbf{h}_k < 0$ , which ensures 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, we decrease the value of the step size

# Summary of descent methods

## General algorithm

- 1: Initialize  $\theta_k$
- 2: **repeat**
- 3:   Find direction of descent  $\mathbf{h}_k \in \mathbb{R}^d$
- 4:   Line search: find the step  $\alpha_k > 0$
- 5:   Update:  $\theta_{k+1} \leftarrow \theta_k + \alpha_k \mathbf{h}_k$
- 6: **until** convergence

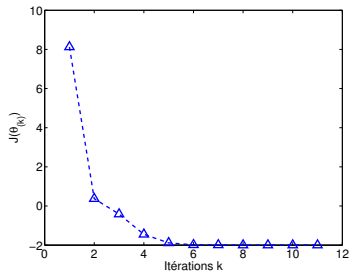
Method	Direction of descent $h$	Complexity	Convergence
Gradient	$-\nabla J(\theta)$	$\mathcal{O}(d)$	linear
Quasi-Newton	$-B(\theta)^{-1} \nabla J(\theta)$	$\mathcal{O}(d^2)$	superlineare
Newton	$-H(\theta)^{-1} \nabla J(\theta)$	$\mathcal{O}(d^3)$	quadratic

- Step size computation: backtracking (common) or optimal step size
- Complexity of each method: depends on the complexity of calculating  $\mathbf{h}$ , the search for  $\alpha$ , and the number of iterations performed till convergence

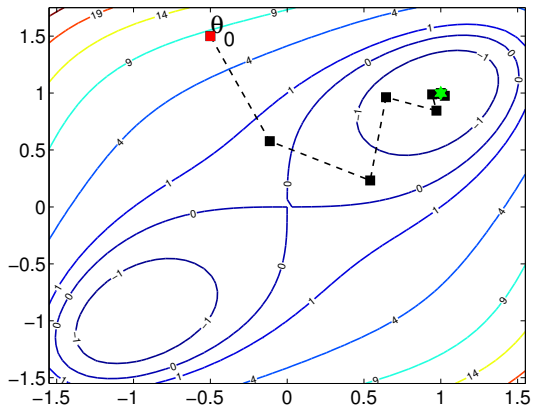


# Gradient method

## J along the iterations

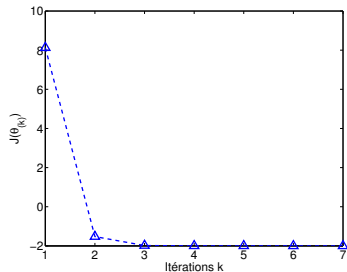


## Evolution of the iterates

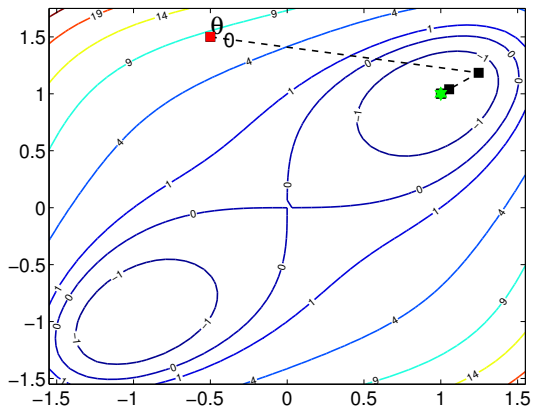


# Newton method

## 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

# 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
- Not covered in this lecture:
  - Convergence analysis of the studied algorithms
  - Non-smooth optimization
  - Gradient-free optimization