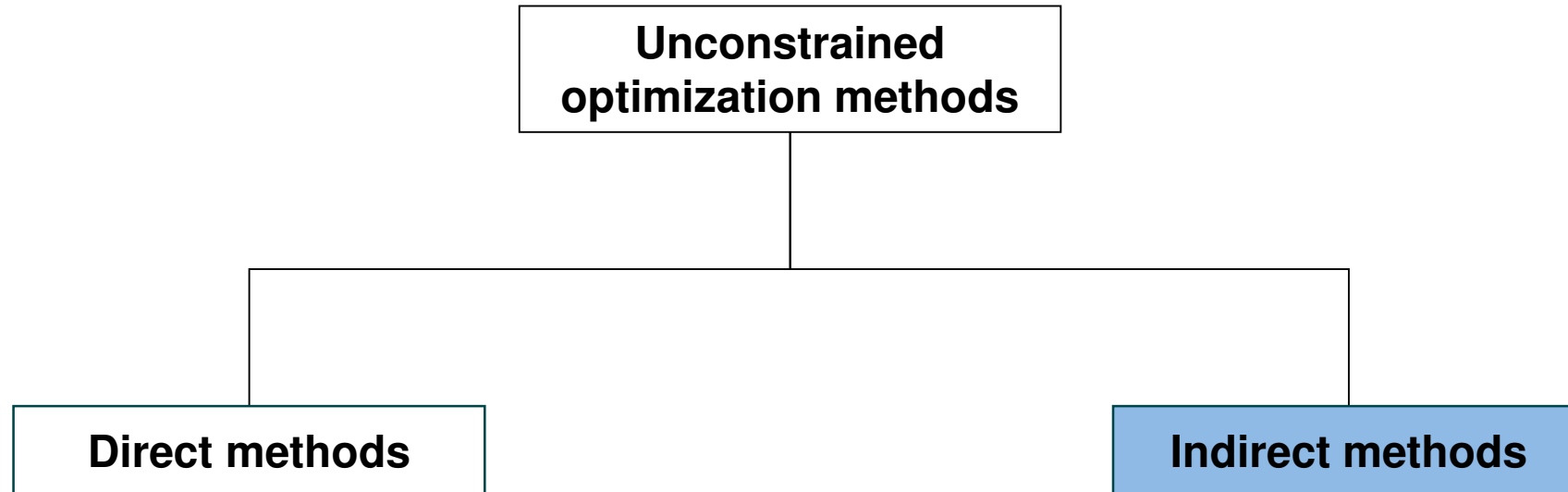




Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

Basic solution methods for unconstrained problems



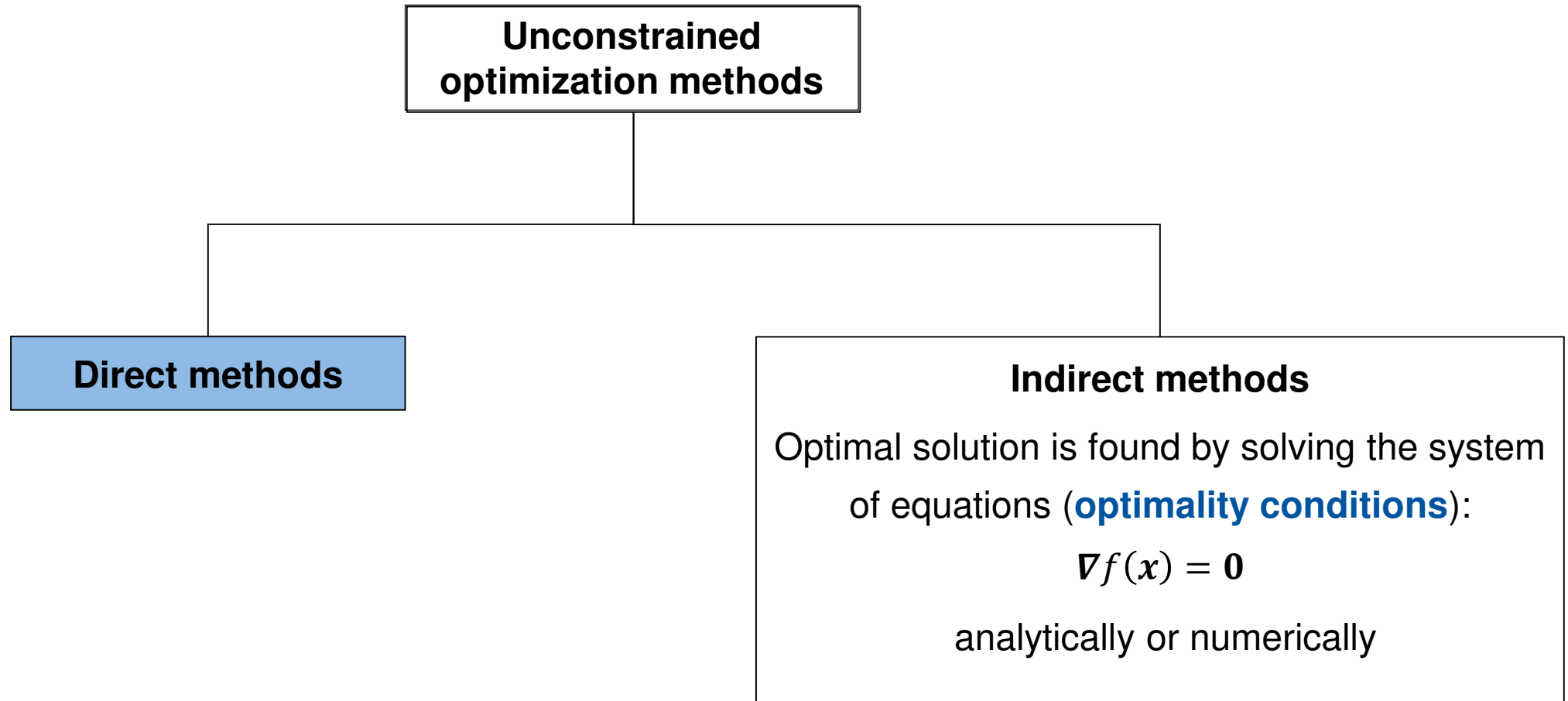
Indirect Methods – Concept

- First-order necessary conditions

$$\nabla f(\mathbf{x}) = \mathbf{0} \Leftrightarrow \begin{cases} \left. \frac{\partial f}{\partial x_1} \right|_{\mathbf{x}} = 0 = g_1(\mathbf{x}) \\ \left. \frac{\partial f}{\partial x_2} \right|_{\mathbf{x}} = 0 = g_2(\mathbf{x}) \\ \vdots \\ \left. \frac{\partial f}{\partial x_n} \right|_{\mathbf{x}} = 0 = g_n(\mathbf{x}) \end{cases} \Leftrightarrow \begin{array}{l} \text{nonlinear system of} \\ \text{equations} \\ \mathbf{g}(\mathbf{x}) = \mathbf{0} \end{array}$$

- The optimal solution is found by [solving the system of equations](#) analytically or numerically (e.g., by Newton's method).
- Differentiation and solution of the system of equations is [challenging for complex problems!](#)

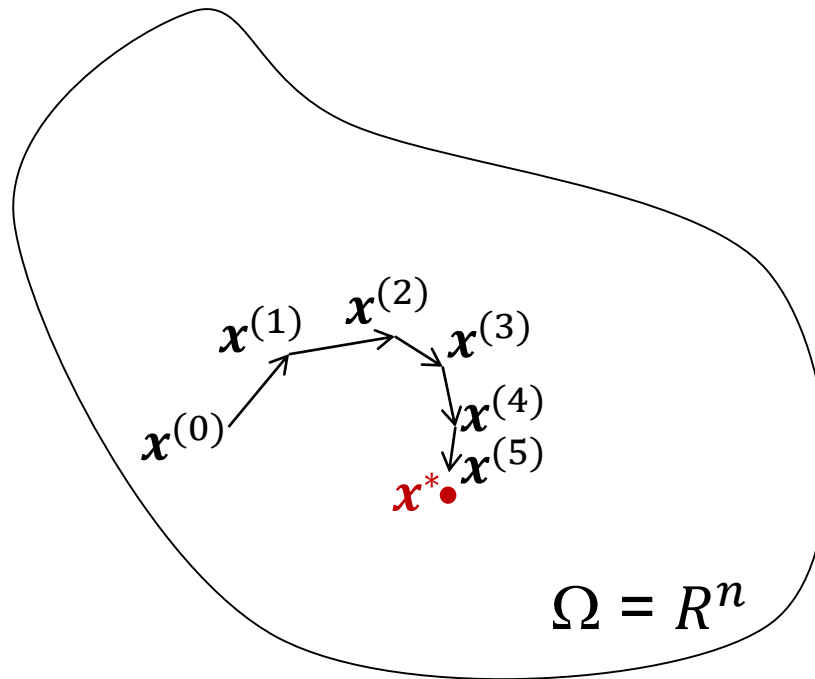
Solution Methods for Unconstrained Optimization



Direct Methods – Concept

Idea: Construct a **convergent** sequence of $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$, which fulfills the following conditions:

$$\exists \bar{k} \geq 0: f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)}) \quad \forall k > \bar{k} \quad \text{and} \quad \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^* \in R^n$$



Definition: Rate of Convergence

Idea: Construct a **convergent** sequence of $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty}$, which fulfills the following conditions:

$$\exists \bar{k} \geq 0: f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)}) \quad \forall k > \bar{k} \quad \text{and} \quad \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^* \in \mathbb{R}^n$$

Rate of convergence:

- **Linear:** if there exists a constant $C \in (0,1)$, such that for sufficiently large k :

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq C \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$$

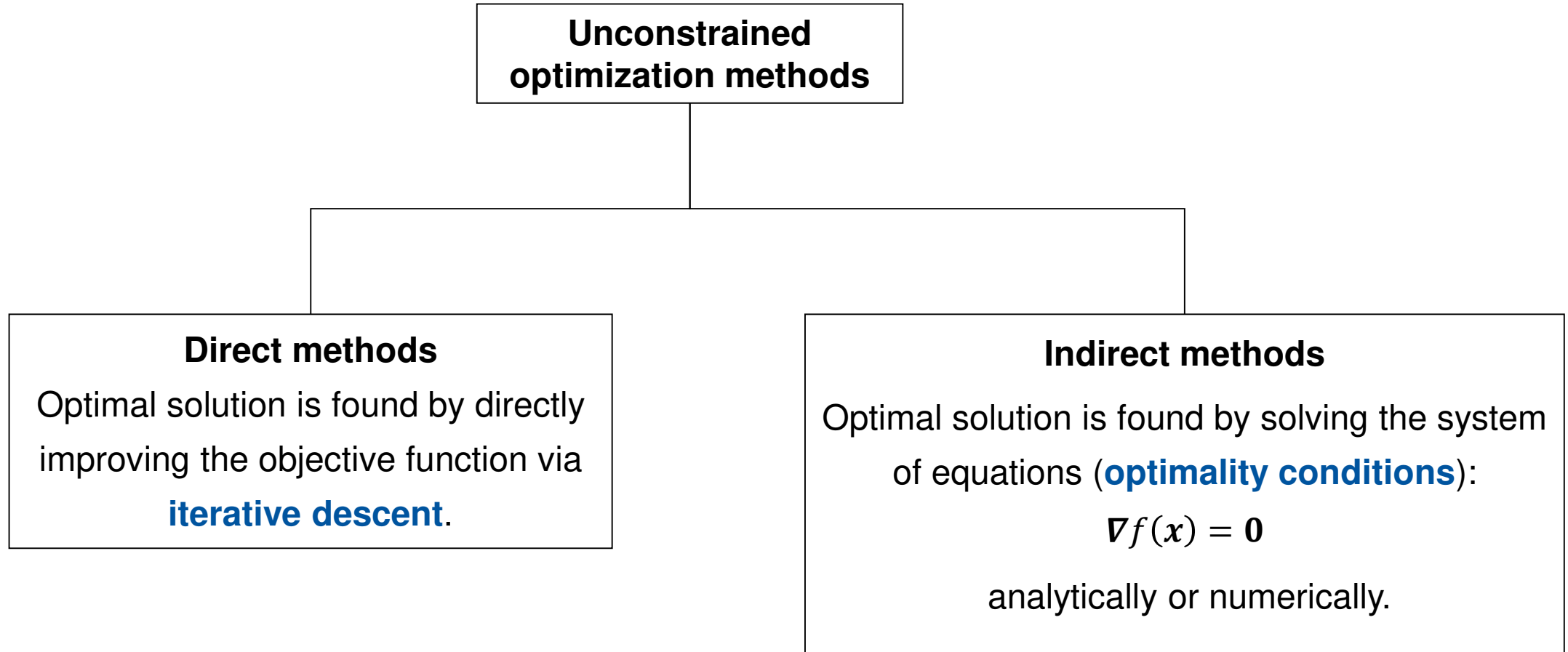
- **Order p** (often $p = 2$): if there exists a constant $M > 0$, such that

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq M \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p$$

- **Superlinear:** if there exists a sequence c_k converging to zero, i.e., $\lim_{k \rightarrow \infty} c_k = 0$, such that

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq c_k \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$$

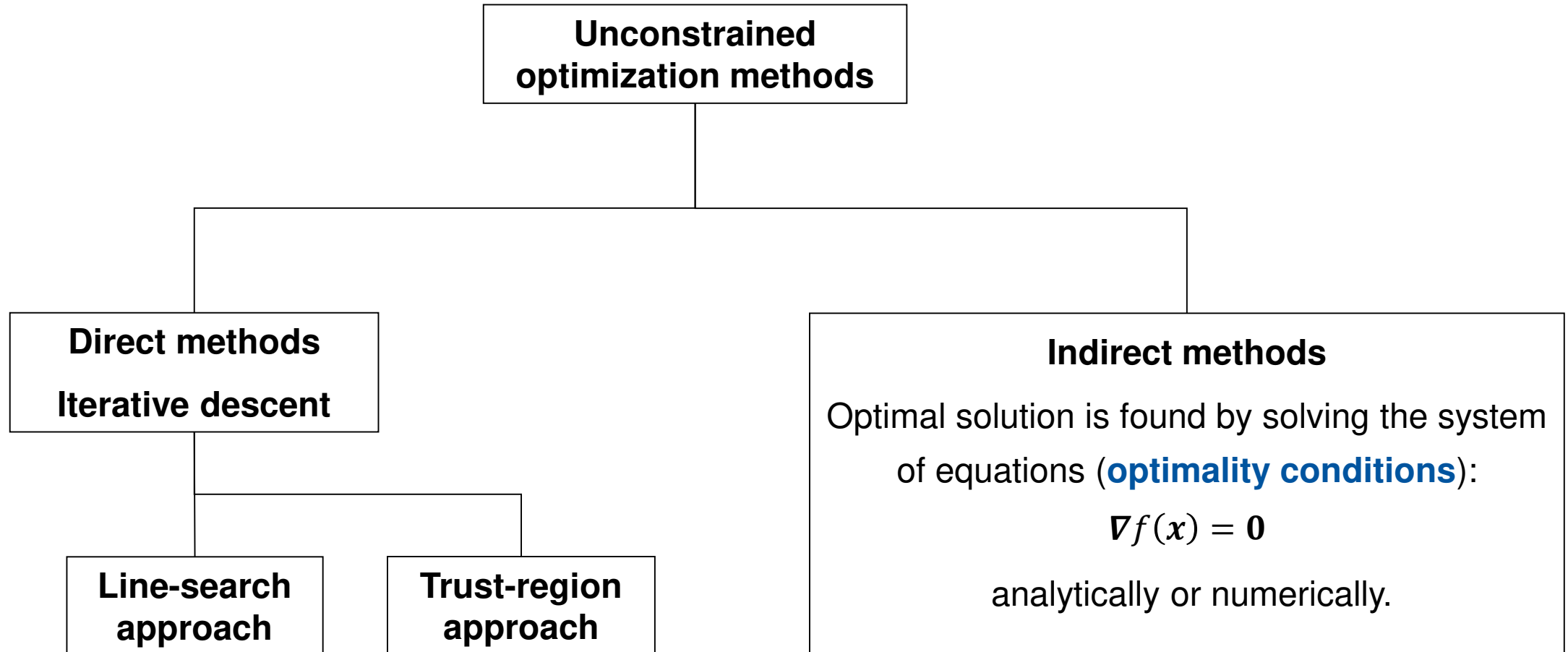
Solution Methods for Unconstrained Optimization



Direct vs Indirect: Nomenclature not consistent in Literature

- Throughout class we use "direct" and "indirect":
 - "indirect methods": 1. set up optimality conditions and then 2. try to solve the system of equations (or equations and inequalities)
 - "direct methods": directly aim to improve objective function (or objective function and constraints). These methods hope to converge to optimality conditions.
- In the literature there are many alternative uses of the word, including
 - exactly the opposite than ours
 - "direct": without the use of derivatives, "indirect": using derivatives
 - only in the context of dynamic optimization problems:
 - "direct": first convert to nonlinear program
 - "indirect": first set up optimality conditions
 - only in the context of constrained problems
 - "direct": only feasible iterates
 - "indirect": infeasible iterates are allowed

Solution Methods for Unconstrained Optimization



Check Yourself

- What are direct vs indirect methods?
- Which direct methods did we learn?
- Which convergence rates exist? Why is the convergence rate important?

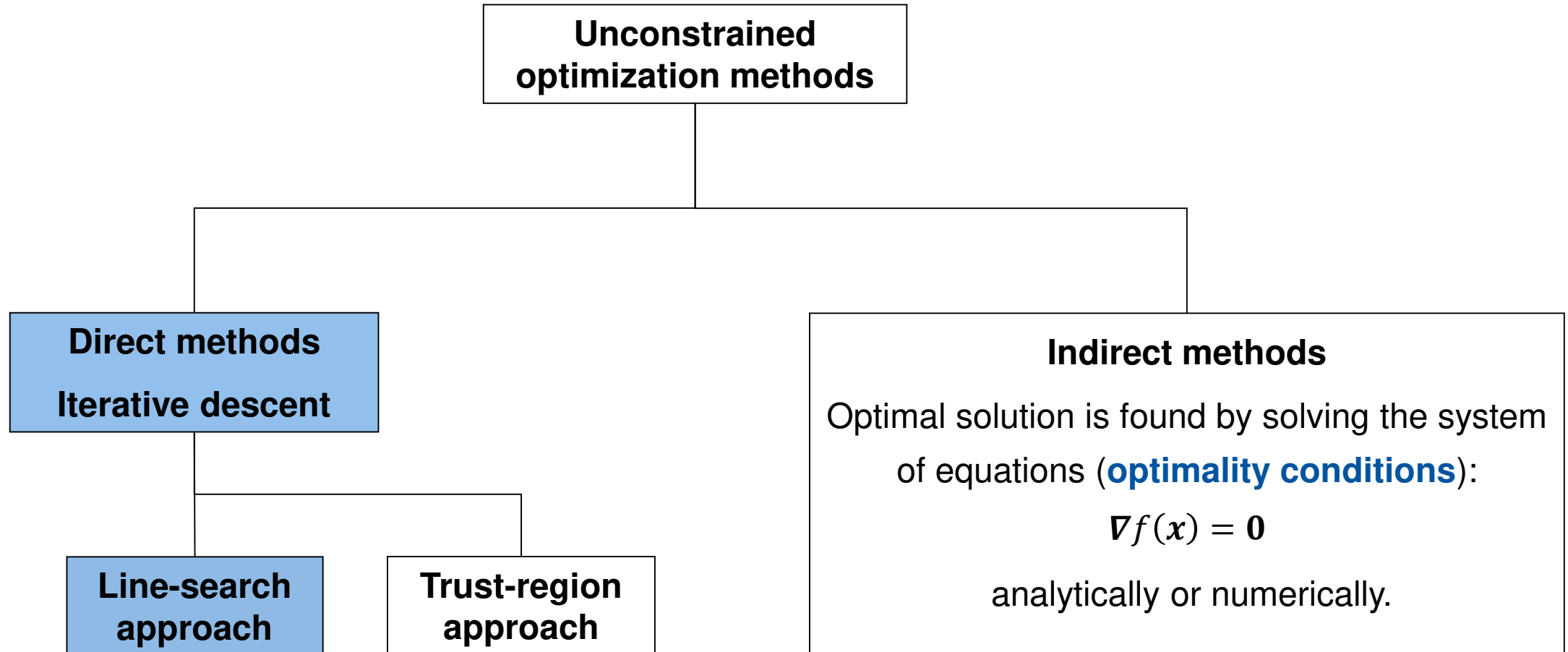


Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

Line search: basic idea and step length

Solution Methods for Unconstrained Optimization



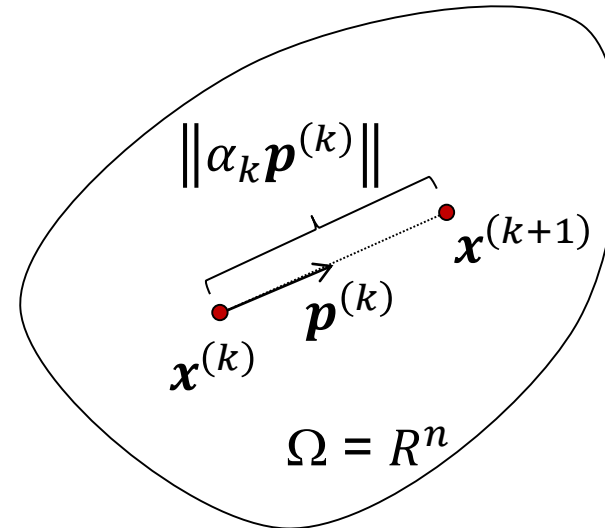
Direct Methods – Line-Search Approach

Definition (descent direction):

A vector \mathbf{p} is called **descent direction** at $\mathbf{x}^{(k)}$, if $\nabla f(\mathbf{x}^{(k)})^T \mathbf{p} < 0$ holds.

Basic algorithm (line-search):

1. **Choose** a descent direction, $\mathbf{p}^{(k)}$, such that
$$\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} < 0$$
2. **Determine** a step length α_k
3. **Set** $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$



Open issues:

- Determination of the descent direction $\mathbf{p}^{(k)}$?
- Calculation of the step length α_k ?

Calculation of Step Length α_k

The exact line search strategy:

1. Define the one-dimensional function along the descent direction $\mathbf{p}^{(k)}$:

$$\phi(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{p}^{(k)})$$

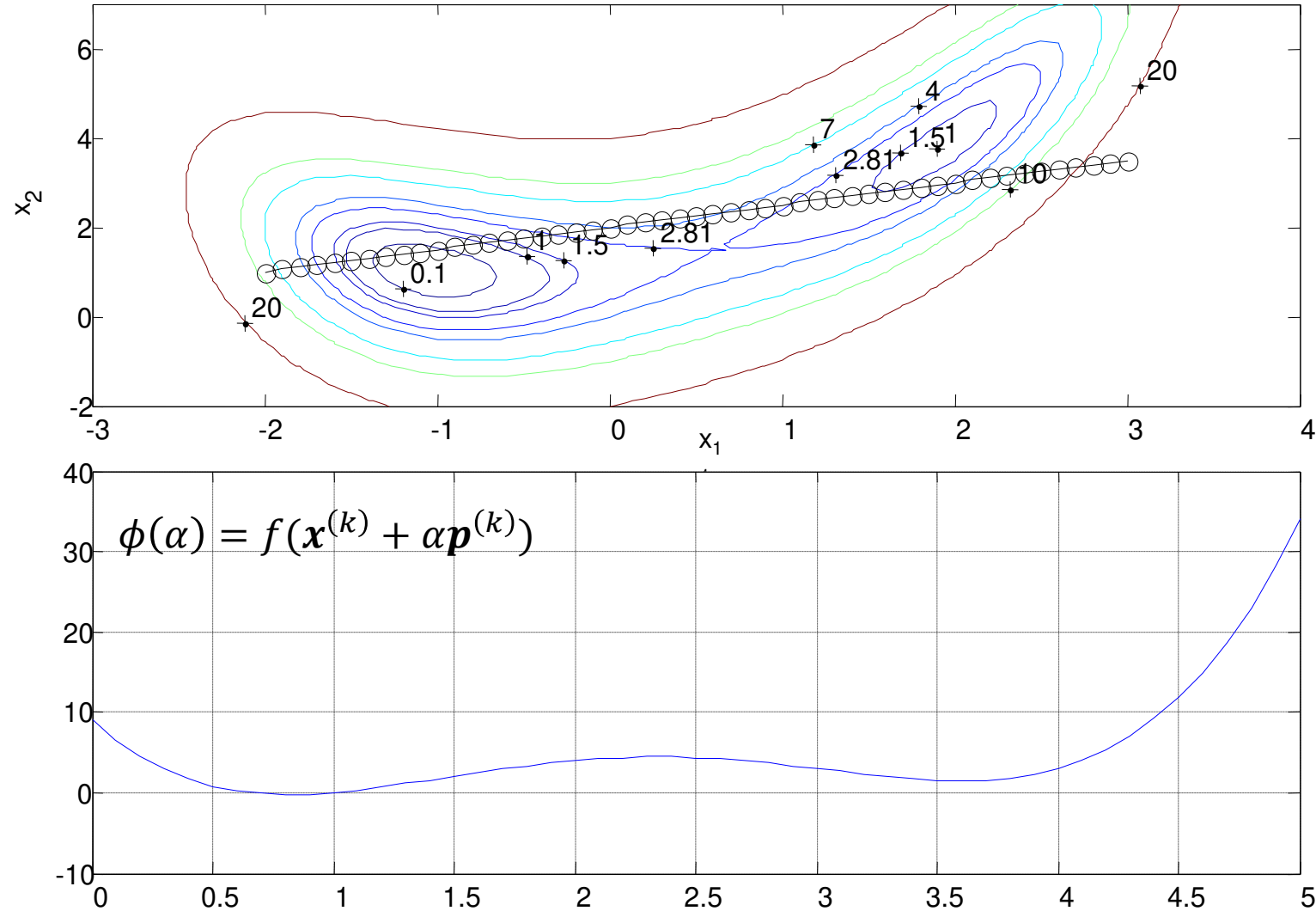
2. Solve the one-dimensional minimization problem

$$\min_{\alpha > 0} \phi(\alpha)$$

Remarks

1. Naively speaking it would be ideal to **globally minimize** $\phi(\alpha)$. Generally, it is very expensive to find this solution. It is not necessarily a good idea since the search is one-dimensional
2. One could also search for some **local solution**. But this is often also too expensive (need function and/or gradient evaluations at a number of points).
3. Practical strategies (so-called **non-exact LS**): find α such that $f(\mathbf{x}^{(k+1)})$ becomes as small as possible with minimal effort.

Practical Line-Search Strategies

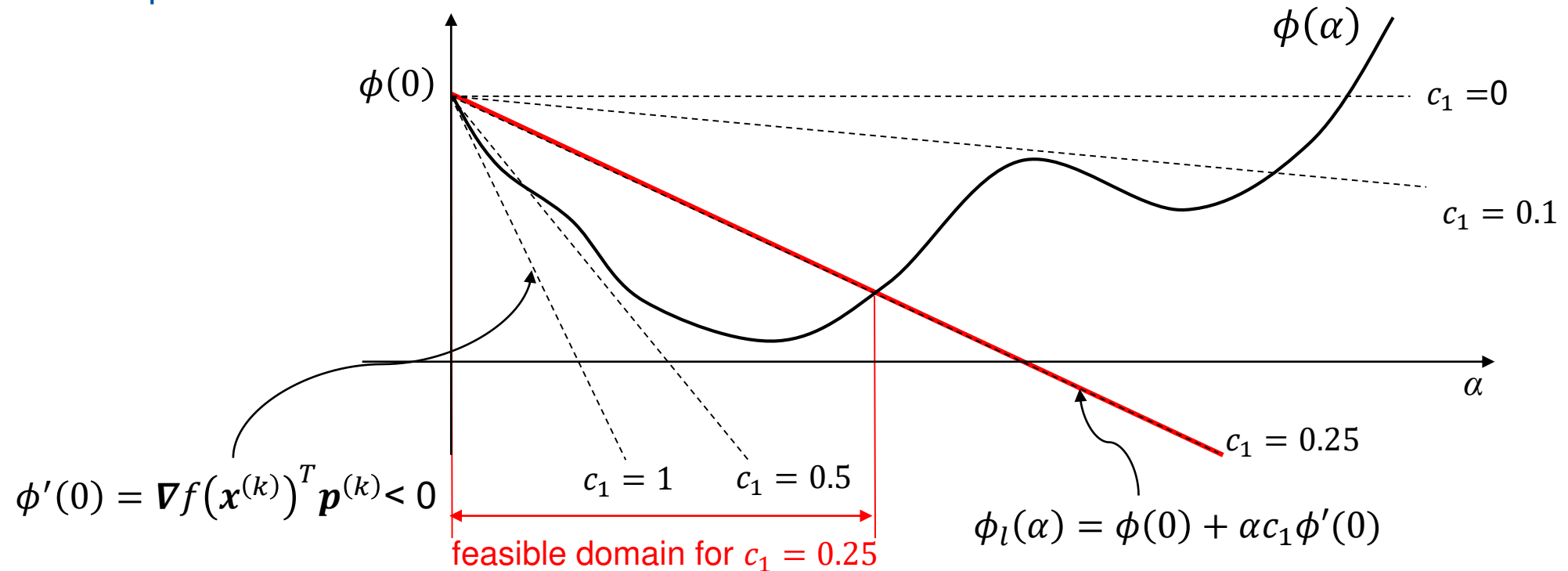


Armijo Condition

Theorem^[1]:

Let f be continuously differentiable, $\mathbf{p}^{(k)}$ a descent direction, and let $c_1 \in (0,1)$ be given. Then there exists an $\alpha > 0$, such that for $\phi(\alpha) := f(\mathbf{x}^{(k)} + \alpha\mathbf{p}^{(k)})$, the condition $\phi(\alpha) \leq \phi(0) + \alpha c_1 \phi'(0)$ holds.

Geometrical interpretation:



Simple Line-Search Algorithm

Remarks:

1. The choice of a step length, which fulfills the Armijo condition guarantees the descent of f :

$$\phi'(0) = \nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} < 0 \quad (\mathbf{p}^{(k)} \text{ is a descent direction})$$

$$\phi(\alpha) \leq \phi(0) + \alpha c_1 \phi'(0)$$

$$\Rightarrow \phi(\alpha) < \phi(0)$$

\Rightarrow a descent is guaranteed!

2. The choice of c_1 is crucial:

- Large c_1 leads to small values of α , such that $\mathbf{x}^{(k+1)} \approx \mathbf{x}^{(k)}$.
- Small c_1 potentially results in small reduction of f and therefore slower convergence

Simple line-search algorithm:

choose $\alpha_1 > 0$; $\rho, c_1 \in (0,1)$

set $\alpha = \alpha_1$

repeat $\alpha \leftarrow \rho \alpha$ until $\phi(\alpha) \leq \phi(0) + \alpha c_1 \phi'(0)$

Improved Line-Search Algorithm

choose $\alpha_0 > 0$ and $c_1 \in (0,1)$

if $\phi(\alpha_0) \leq \phi(0) + \alpha_0 c_1 \phi'(0)$ STOP, else

find a better $\alpha \in (0, \alpha_0)$ through *quadratic interpolation* of available data:

$$\alpha_1 = -\frac{\phi'(0)\alpha_0^2}{2[\phi(\alpha_0) - \phi(0) - \phi'(0)\alpha_0]}$$

if $\phi(\alpha_1) \leq \phi(0) + \alpha_1 c_1 \phi'(0)$ STOP, else

find a better $\alpha \in (0, \alpha_1)$ through *cubic interpolation* of available data (how ?)

repeat the procedure of *cubic interpolation*, until the condition is fulfilled

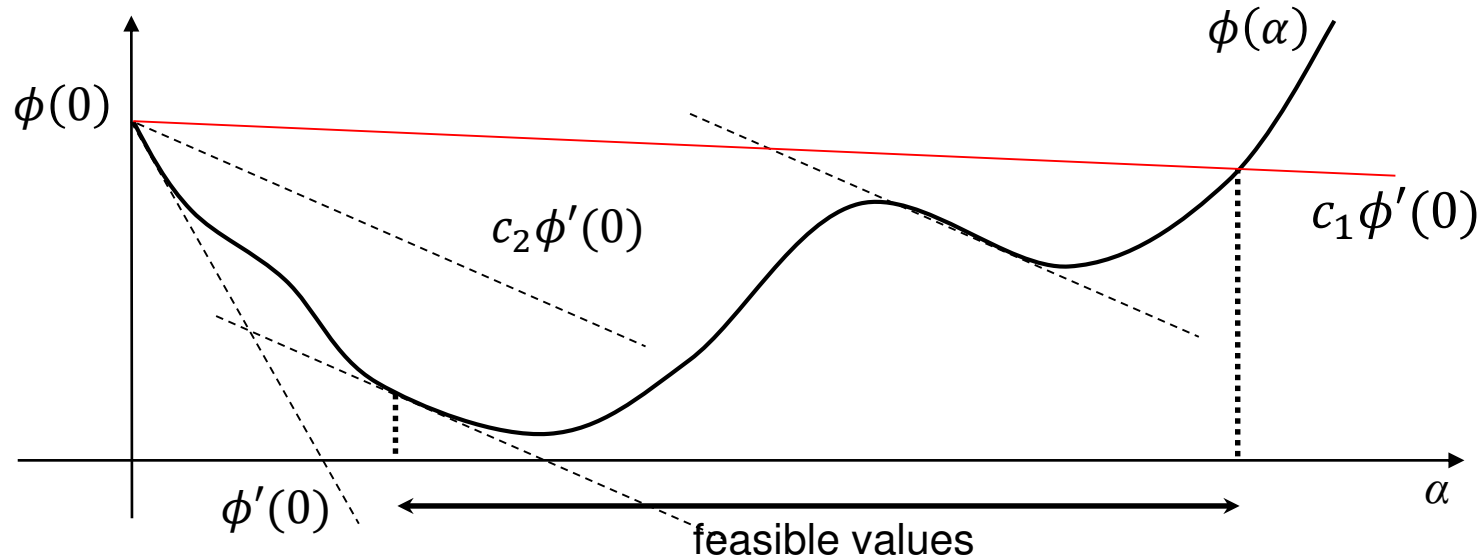
Wolfe Conditions

Theorem^[1]:

Let f be continuously differentiable, $\mathbf{p}^{(k)}$ a descent direction and $c_1 \in (0,1)$, $c_2 \in (c_1, 1)$. Then, there exists an $\alpha > 0$, such that $\phi(\alpha) \leq \phi(0) + \alpha c_1 \phi'(0)$

$$\phi'(\alpha) \geq c_2 \phi'(0) \text{ (slope condition)}$$

Geometric interpretation: → guarantee minimum step length!



Relevance:

Wolfe Conditions promote convergence to a stationary point^[1]

Check Yourself

- Explain the basic ideas of the line-search method.
- What is a descent direction? How it is defined?
- Explain the Armijo-rule and its potential drawbacks?
- Explain the Wolfe conditions and the advantage compared to Armijo's rule.

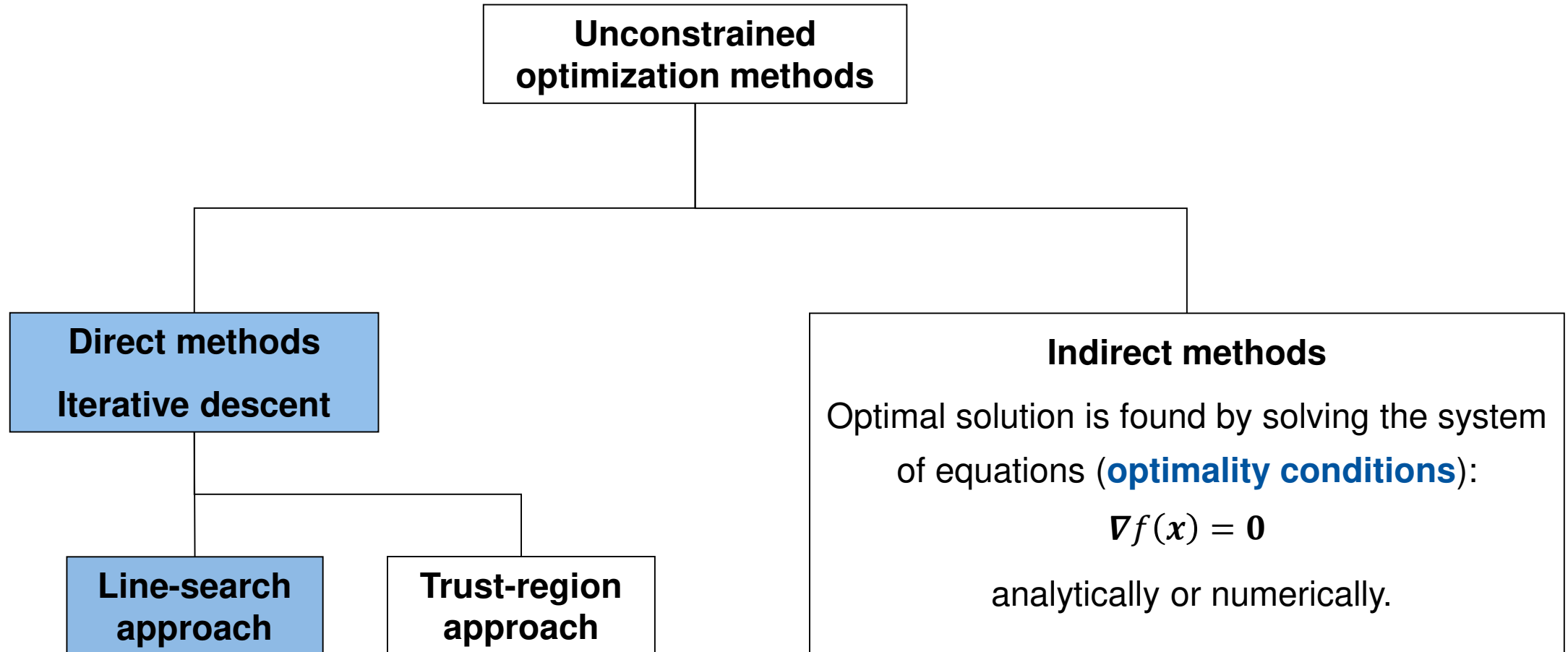


Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

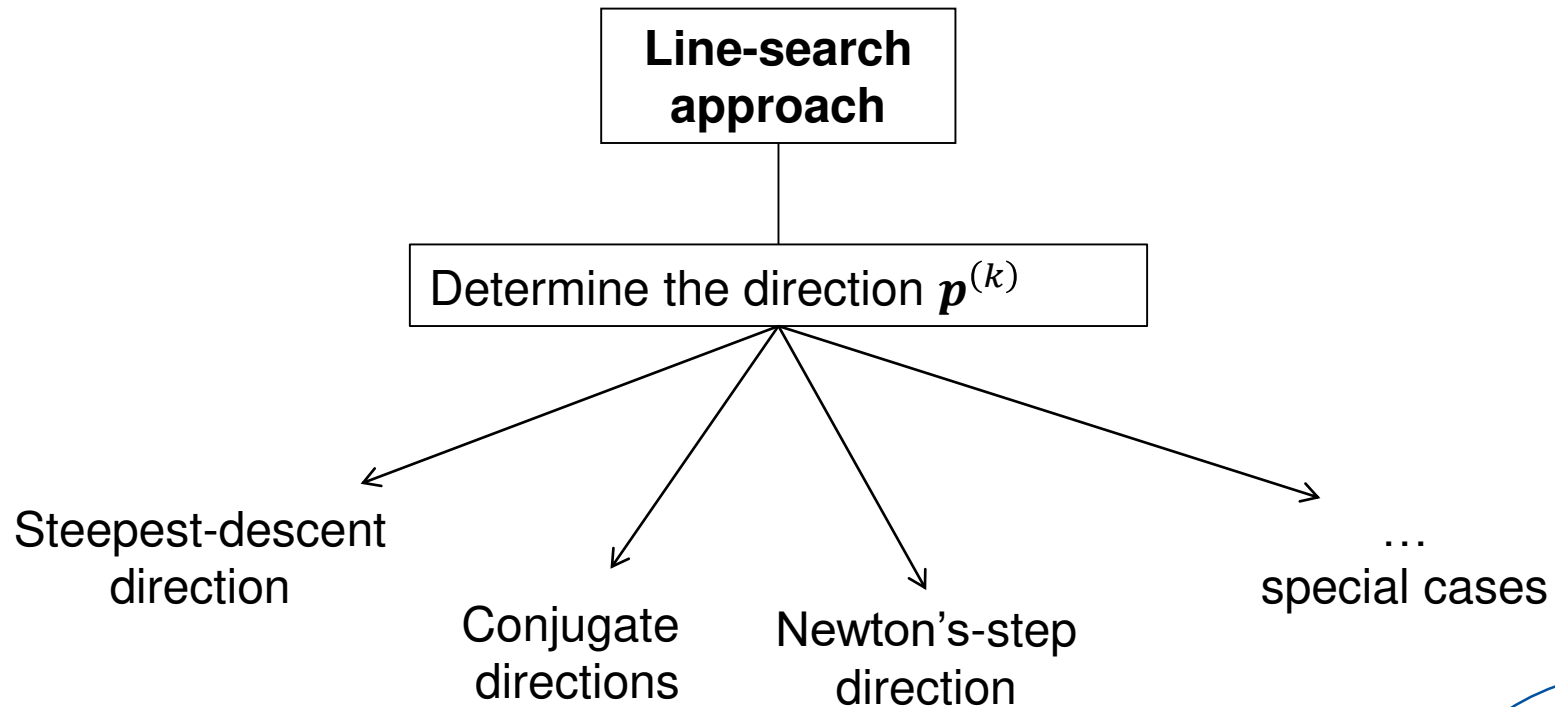
Line search: simple directions

Solution Methods for Unconstrained Optimization



Determination of a Descent Direction: A Toolbox

Line-search approaches differ from each other with respect to the determination of descent direction and step length.



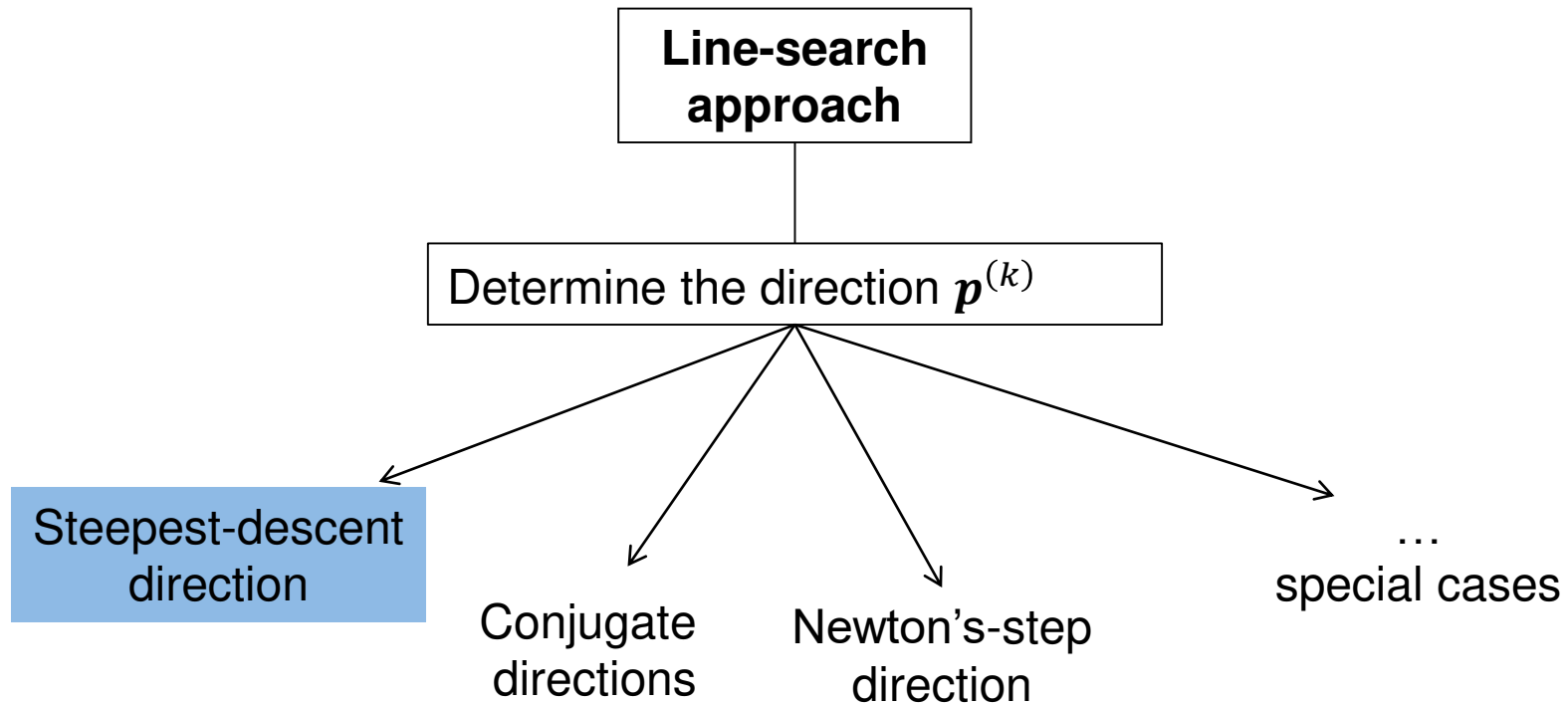
Many gradient methods use a symmetric positive definite matrix $\mathbf{D}^{(k)}$ and calculate

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{D}^{(k)} \nabla f(\mathbf{x}^{(k)})$$

*Extra work:
prove that it
guarantees
descent!*

Determination of a Descent Direction: A Toolbox

Line-search approaches differ from each other with respect to the determination of descent direction and step length.



Many gradient methods use a symmetric positive definite matrix $\mathbf{D}^{(k)}$ and calculate

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{D}^{(k)} \nabla f(\mathbf{x}^{(k)})$$

Steepest-Descent Direction (1)

Taylor series: $f(\mathbf{x}^{(k)} + \alpha \mathbf{p}^{(k)}) = f(\mathbf{x}^{(k)}) + \boxed{\alpha \nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}} + O(\alpha^2)$

The rate of change of f at $\mathbf{x}^{(k)}$ along the direction $\mathbf{p}^{(k)}$ is the coefficient in the linear term:

$$\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}$$

The unit direction $\mathbf{p}^{(k)}$ with the **highest rate of change** is the solution of the following problem

$$\min_{\mathbf{p}^{(k)} \in \mathbb{R}^n} \nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} \quad \text{s. t. } \|\mathbf{p}^{(k)}\| = 1$$

Note that $\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} = \|\nabla f(\mathbf{x}^{(k)})\| \|\mathbf{p}^{(k)}\| \cos(\theta)$

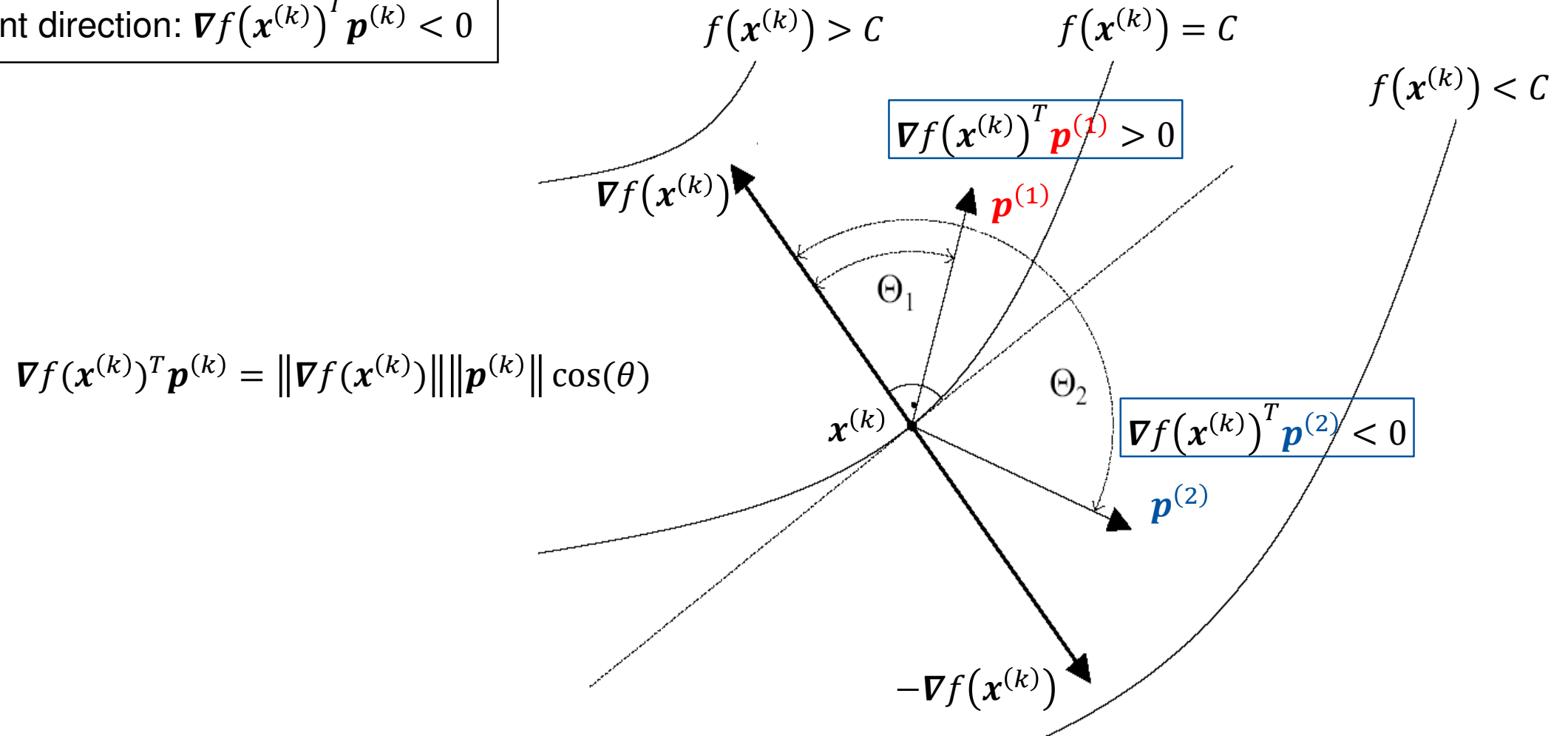
The solution of the problem is achieved for $\cos(\theta) = -1 \Rightarrow \theta = \pi$

$$\Rightarrow \mathbf{p}^{(k)} = -\nabla f(\mathbf{x}^{(k)}) / \|\nabla f(\mathbf{x}^{(k)})\|$$

The choice of $\mathbf{D}^{(k)}$ is the identity matrix \mathbf{I} .

Steepest-Descent Direction (2)

descent direction: $\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} < 0$



Method of Steepest-Descent

Algorithm:

choose $\mathbf{x}^{(0)}$

for $k=0,1,\dots$

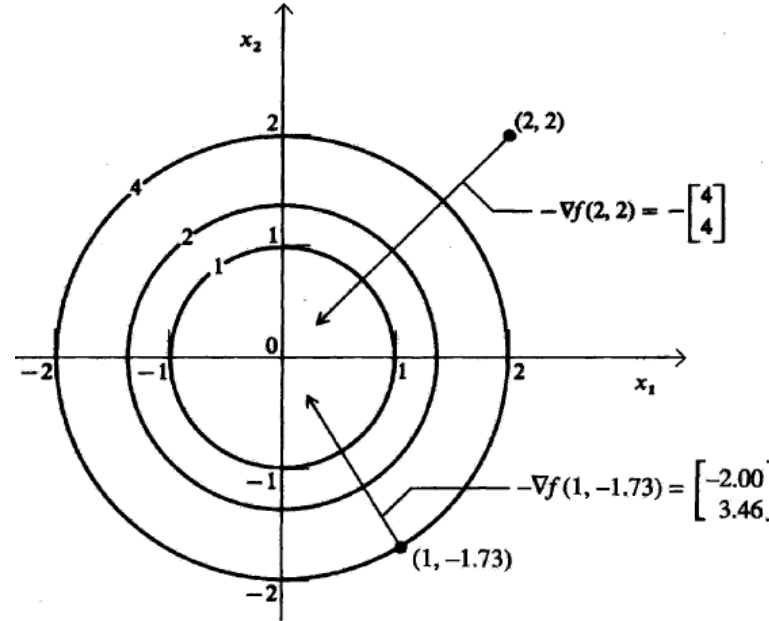
if $\|\nabla f(\mathbf{x}^{(k)})\| \leq \varepsilon$ stop, else

set $\mathbf{p}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$

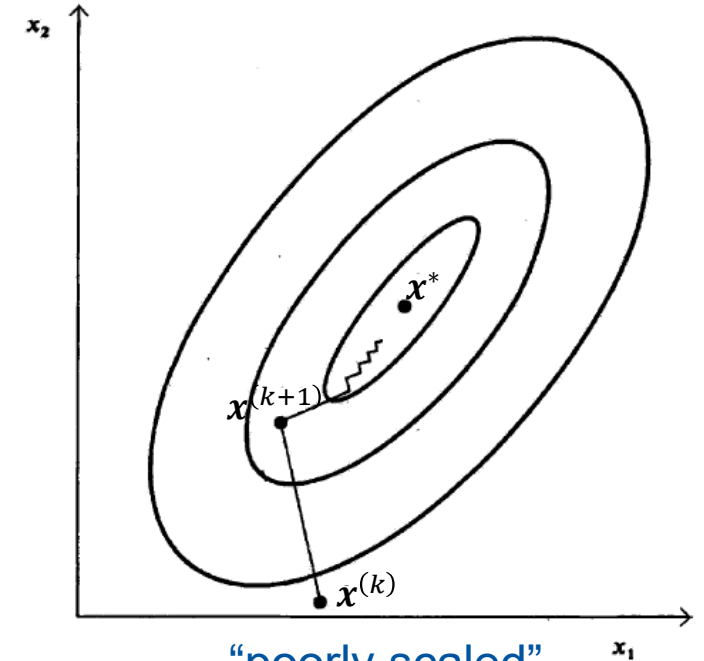
determine the step length α_k (e.g.
using the Armijo rule)

set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$

end for



“well scaled”

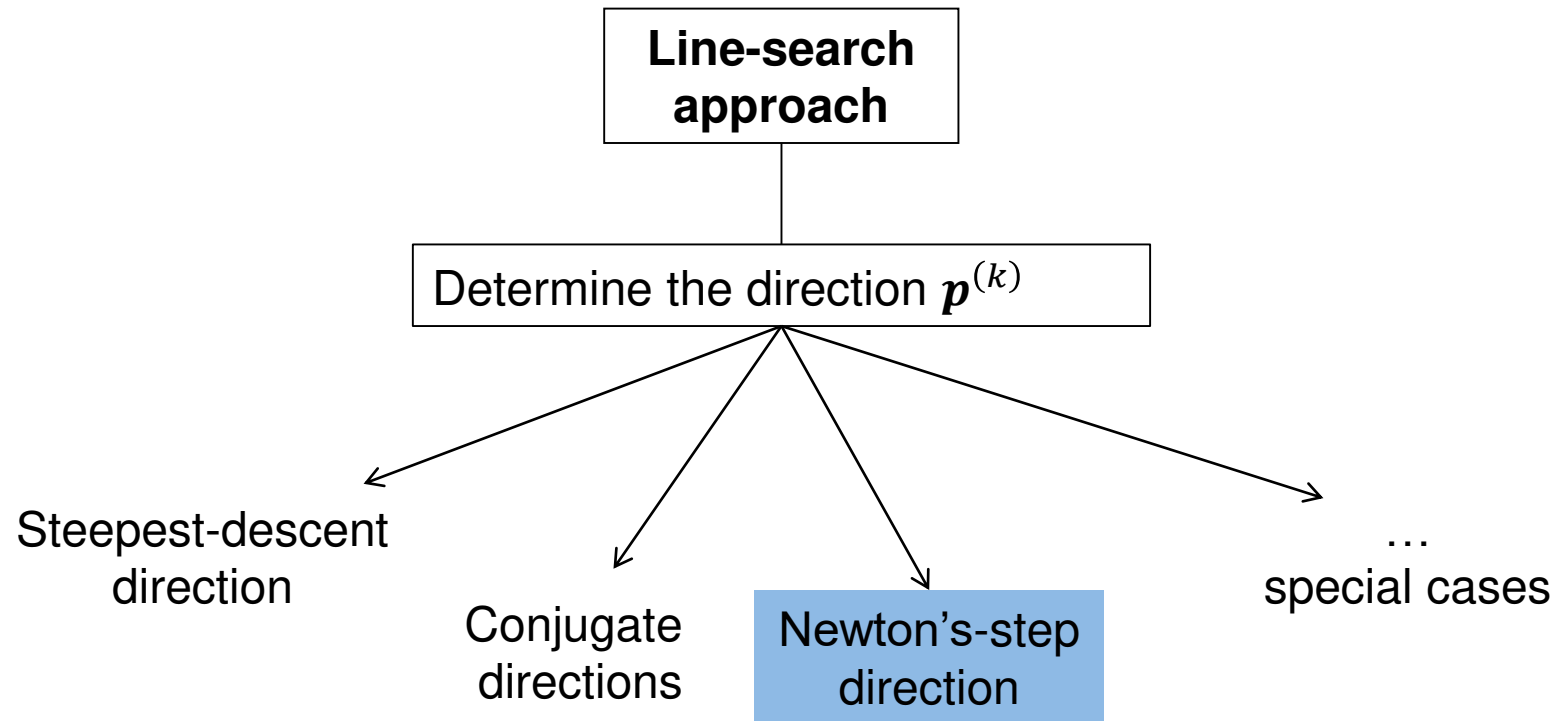


“poorly scaled”

Directions become perpendicular

Determination of a Descent Direction: A Toolbox

Line-search approaches differ from each other with respect to the determination of descent direction and step length.



Many gradient methods use a symmetric positive definite matrix $\mathbf{D}^{(k)}$ and calculate

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{D}^{(k)} \nabla f(\mathbf{x}^{(k)})$$

Newton's Descent Direction

Quadratic approximation of f at $\mathbf{x}^{(k+1)}$

$$m(\mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + \frac{1}{2} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^T \nabla^2 f(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

(1st nec. opt. cond. for m)



$$0 = \nabla m(\mathbf{x}^{(k+1)}) = \nabla f(\mathbf{x}^{(k)}) + \nabla^2 f(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

$$\Rightarrow \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - [\nabla^2 f(\mathbf{x}^{(k)})]^{-1} \nabla f(\mathbf{x}^{(k)})$$



$$\Rightarrow \mathbf{p}^{(k)} = -[\nabla^2 f(\mathbf{x}^{(k)})]^{-1} \nabla f(\mathbf{x}^{(k)})$$

$$\Rightarrow \alpha_k = 1$$

The choice of $\mathbf{D}^{(k)}$ is the inverse of the Hessian

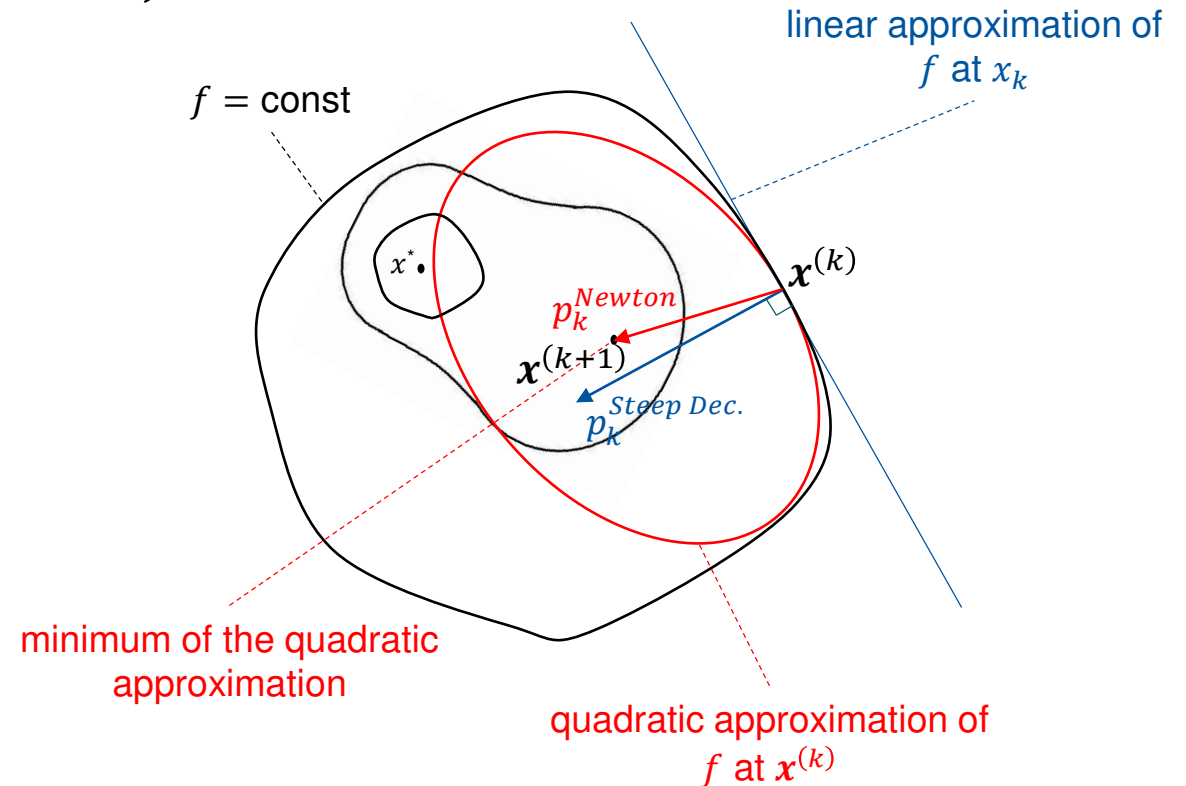


Fig.: Comparison of steepest-descent with Newton's method from viewpoint of objective function approximation

Newton's Method

Algorithm:

choose $\mathbf{x}^{(0)}$

for $k=0,1,\dots$

if $\|\nabla f(\mathbf{x}^{(k)})\| \leq \varepsilon$ stop, else

set $\mathbf{p}^{(k)} = -[\nabla^2 f(\mathbf{x}^{(k)})]^{-1} \nabla f(\mathbf{x}^{(k)})$

set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^{(k)}$

end for

Remarks:

1. line-search?

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$$

$$\mathbf{p}^{(k)} = -[\nabla^2 f(\mathbf{x}^{(k)})]^{-1} \nabla f(\mathbf{x}^{(k)})$$

$$\alpha_k = 1$$

2. (+) locally **quadratic convergence**, if $\mathbf{x}^{(k)}$ close to \mathbf{x}^*
(−) 2nd derivatives & inversion (expensive for large system of equations)
3. If f is quadratic, the algorithm converges in **one iteration**.
4. Convergence to a minimum is **not guaranteed**! Why?

Check Yourself

- Explain the basic ideas of the line-search method.
- Explain the steepest descent method.
- What additional requirements puts Newton's method on the objective function?
- Explain the Newton direction. Is it better than other descent directions? Why is the Newton step-length equal to one?

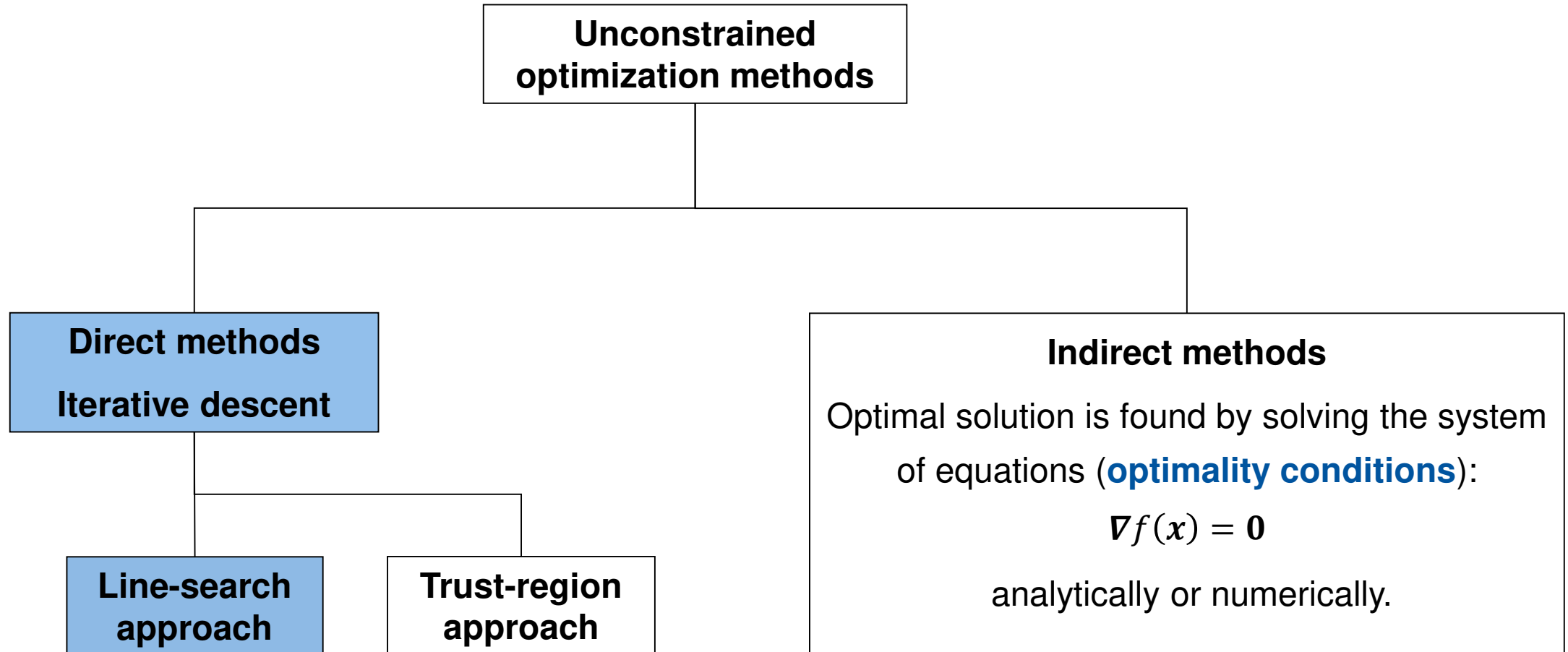


Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

Line search: complexity and examples

Solution Methods for Unconstrained Optimization



Complexity Analysis

Nesterov (2004) proves: “In **general**, optimization problems are **unsolvable**” *

Let F denote a class of problems, e.g., Lipschitz-continuous functions with Lipschitz-constant L , i.e., $|f(x) - f(y)| < L\|x - y\|$, L is assumed to be fixed for all $P \in F$.

“Performance of a method M on a problem $P \in F$ is the total amount of computational effort that is required by M to solve P .” *

“To solve the problem means to find an approximate solution to P with an accuracy $\varepsilon > 0$.” *

For unconstrained problems, the accuracy $\varepsilon > 0$ can be defined as the **norm of the objective’s gradient**.

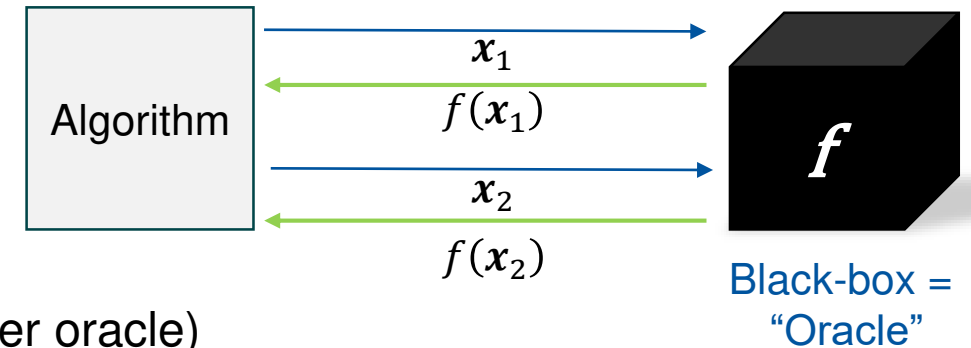
* Yurii Nesterov, *Introductory Lectures on Convex Optimization – A Basic Course*, Kluwer Academic Publishers, (2004)

Complexity Analysis – Measuring Computational Effort

Unit of measurement: Query to an oracle

It is assumed that the objective function is unknown and that the algorithm solves the optimization problem by *querying an oracle* for local information about the unknown objective function. An oracle is simply a “black box” capable of answering any query of the form:

- Given x return the value $f(x)$ (Zeroth-order oracle)
- Given x return $f(x)$ and gradient $\nabla f(x)$ (First-order oracle)
- Given x return $f(x)$, $\nabla f(x)$ and Hessian $\nabla^2 f(x)$ (Second-order oracle)



Analytical Complexity: The smallest number of queries to an oracle to solve Problem P to accuracy ε .^[1]

Arithmetical Complexity: The smallest number of arithmetic operations (including work of the oracle and work of method), required to solve problem P up to accuracy ε .^[1]

Analytical Complexity of Steepest Descent Method

Algorithm:

choose $\mathbf{x}^{(0)}$

for $k=0, 1, \dots$

if $\|\nabla f(\mathbf{x}^{(k)})\| \leq \varepsilon$ stop, else

set $\mathbf{p}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$

determine the step length α_k (e.g. using the Armijo rule)

set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$

end for

- **Problem class:** f is continuously differentiable and $\nabla f(\mathbf{x})$ is Lipschitz-continuous with fixed Lipschitz constant L , i.e., $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| < L\|\mathbf{x} - \mathbf{y}\|$
- **First-order oracle:** returns $f(\mathbf{x})$ and gradient $\nabla f(\mathbf{x})$
- **Worst-case analytical complexity (queries to oracle):** $O\left(\frac{1}{\varepsilon^2}\right)$

Rosenbrock Function

$$\min_{x \in \mathbb{R}^2} f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

- Solution point is $x^* = (1, 1)^T$ - why?

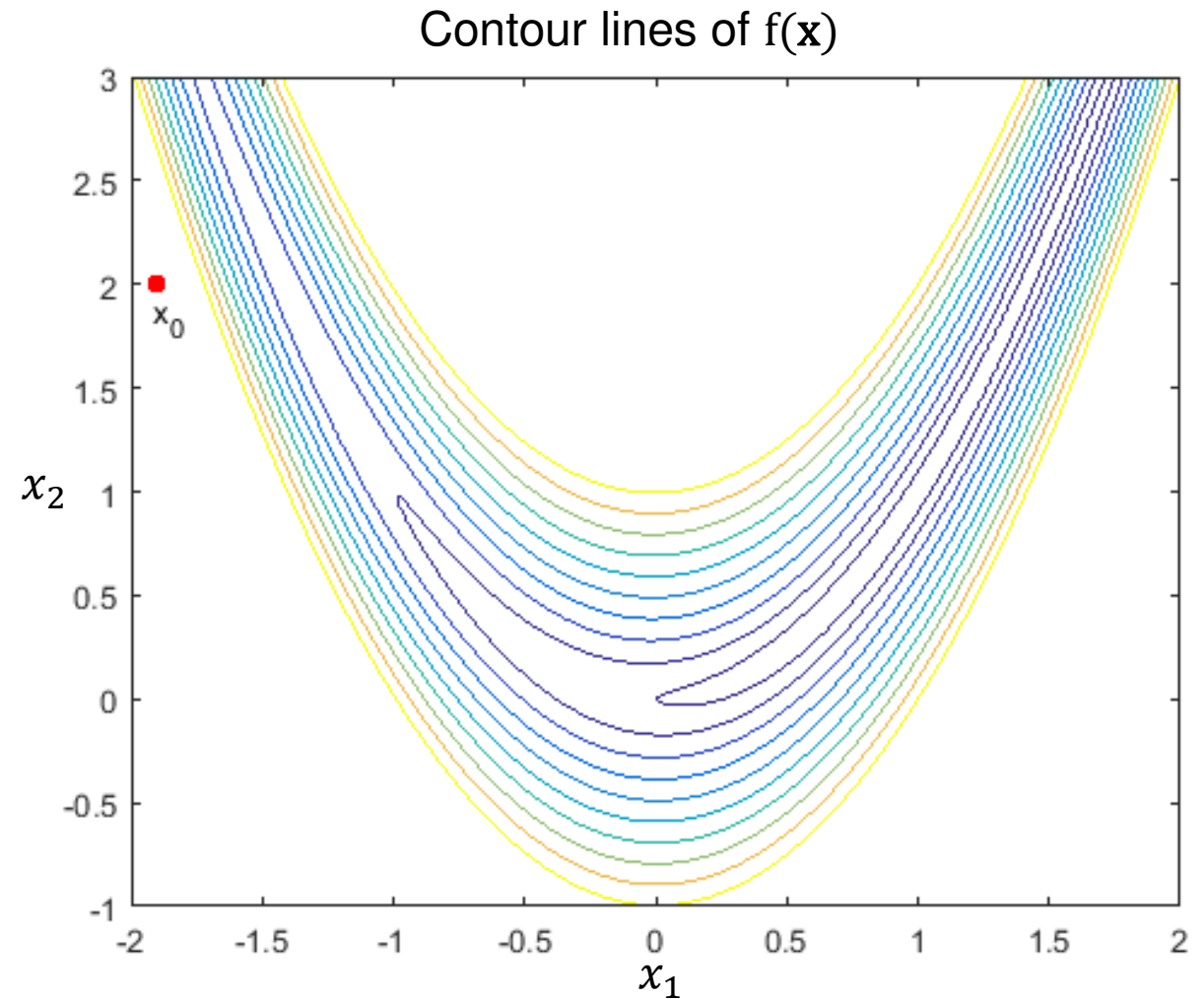
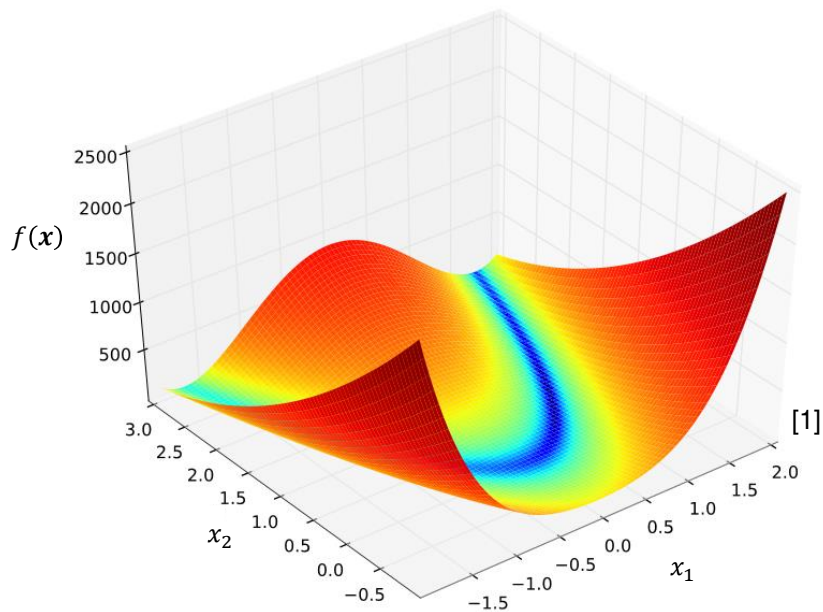


Illustration of Convergence (1)

$$\min_{x \in \mathbb{R}^2} f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

- Steepest descent with **Armijo** line-search

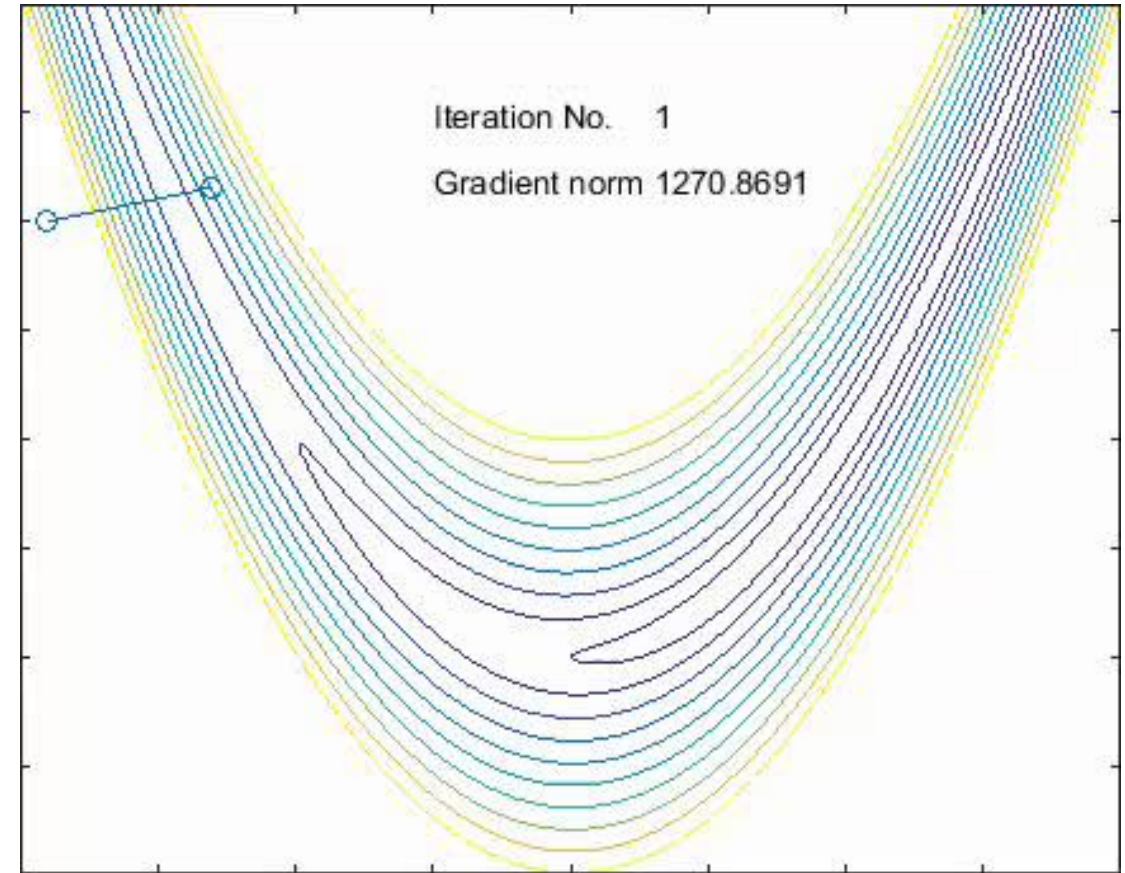
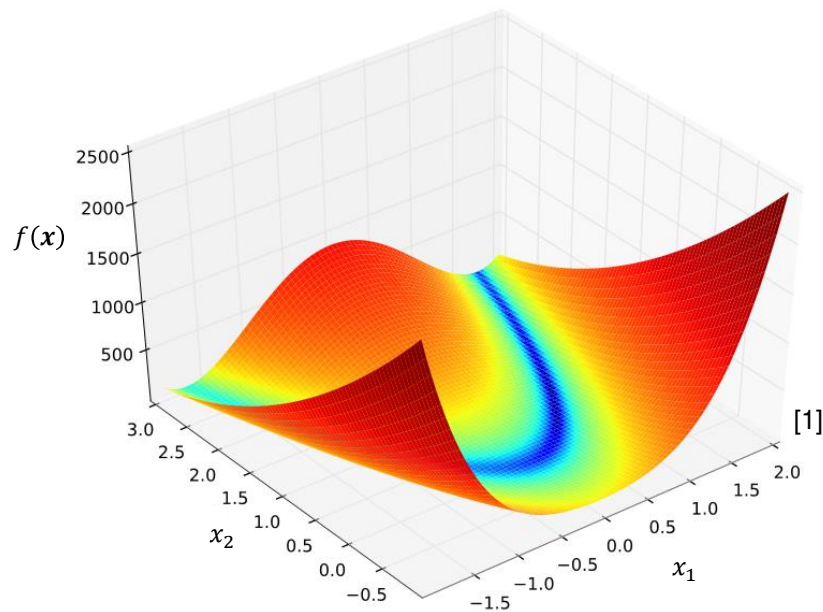
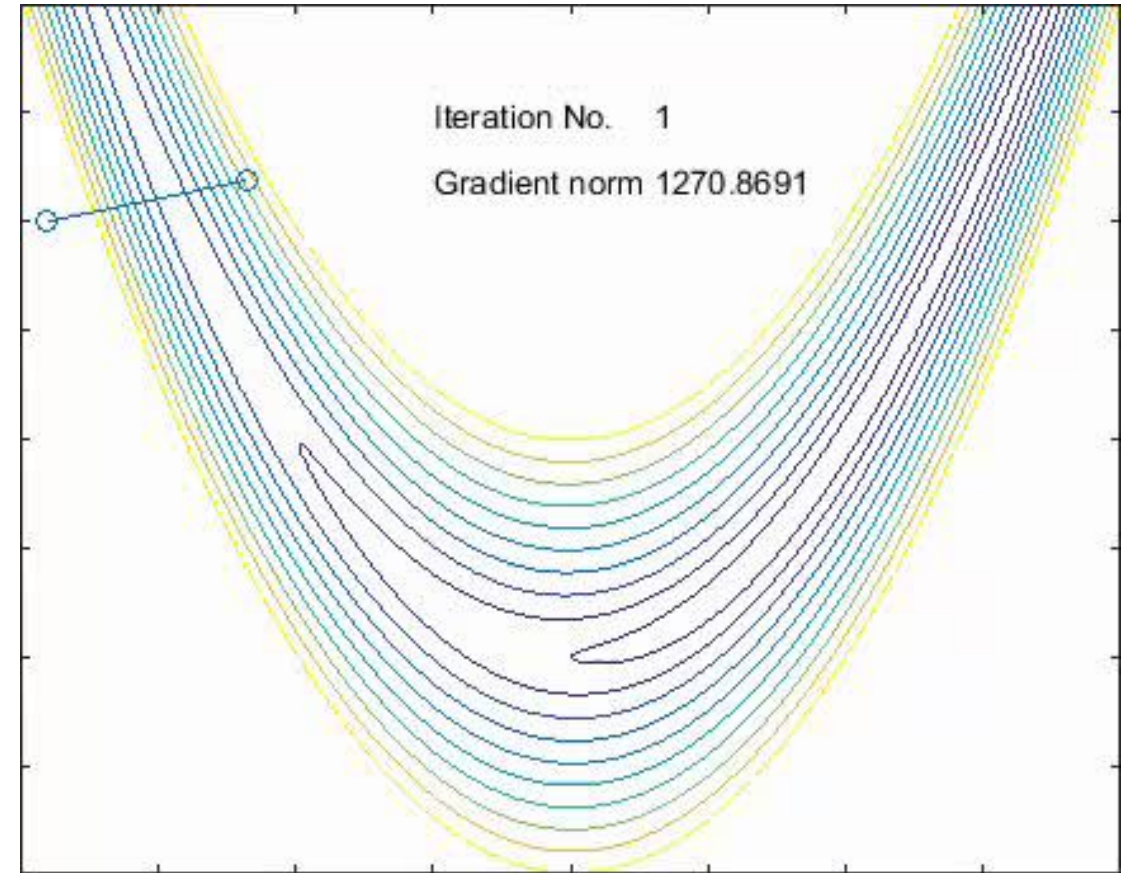
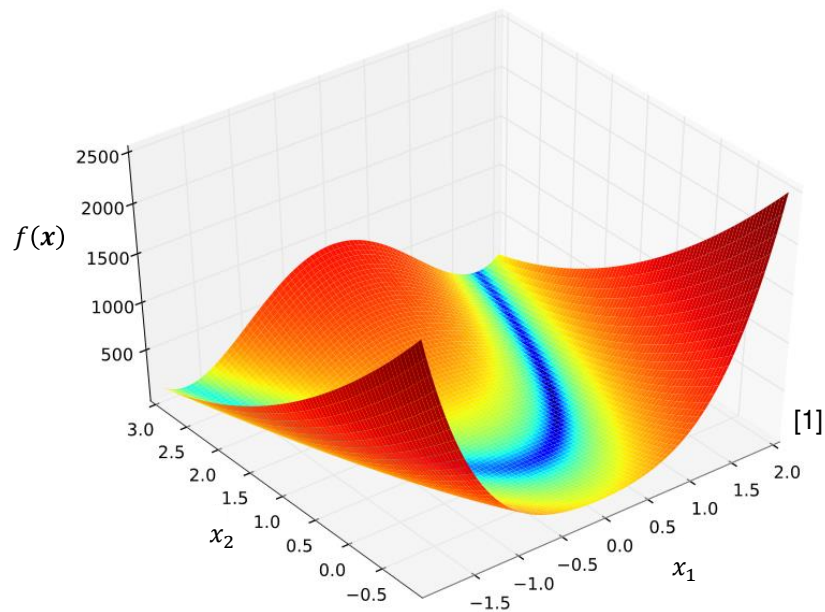


Illustration of Convergence (2)

$$\min_{x \in \mathbb{R}^2} f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

- Steepest descent with **Wolfe** line-search



Analytical Complexity of Newton's Method

- **Problem class:** f is twice continuously differentiable and $\nabla^2 f(\mathbf{x})$ is Lipschitz-continuous with fixed Lipschitz constant L , i.e., $\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| < L\|\mathbf{x} - \mathbf{y}\|$
- **Second-order oracle:** returns $f(\mathbf{x})$, $\nabla f(\mathbf{x})$ and Hessian $\nabla^2 f(\mathbf{x})$

- **Quadratic approximation** of f around $\mathbf{x}^{(k)}$, **line search**

$$m(\mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + \frac{1}{2} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^T \nabla^2 f(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

- **Worst-case analytical complexity:** $O\left(\frac{1}{\varepsilon^{2-\tau}}\right)$, $1 > \tau > 0$, arbitrary but fixed for a given problem

- **Quadratic approximation** of f around $\mathbf{x}^{(k)}$ with **cubic regularization**, **line search**

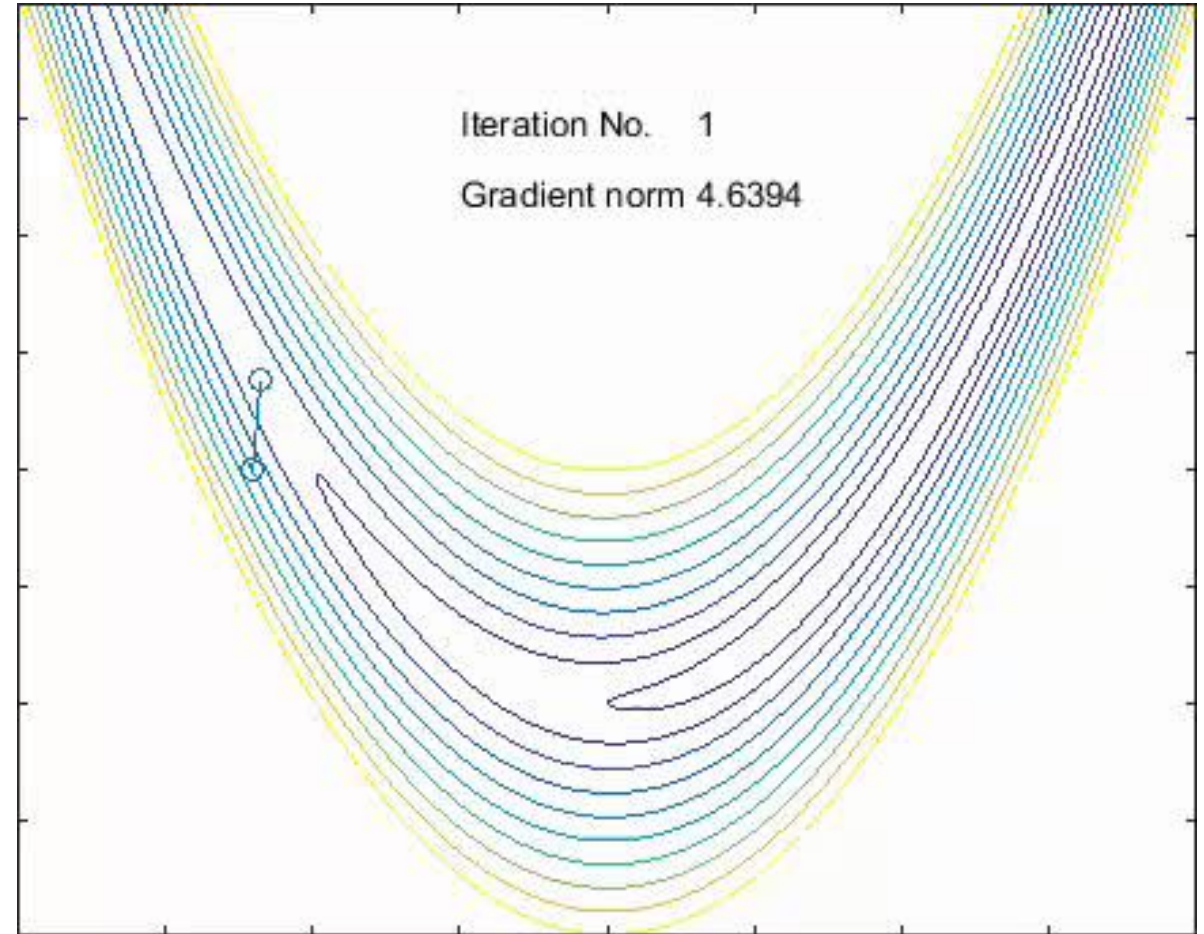
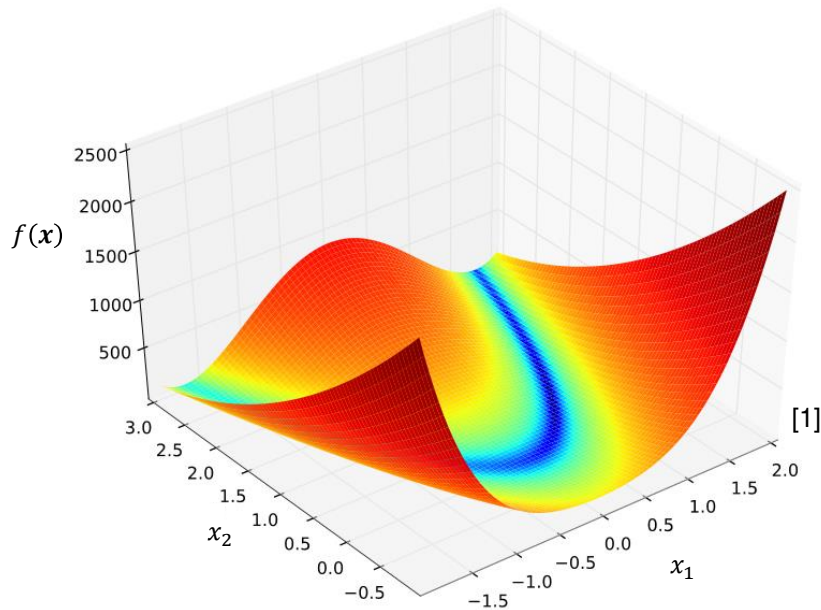
$$m_{\text{regularized}}(\mathbf{x}^{(k+1)}) = m(\mathbf{x}^{(k+1)}) + \frac{1}{3} \sigma_k \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|^3$$

- **Worst-case analytical complexity:** $O\left(\frac{1}{\varepsilon^{3/2}}\right)$

Illustration of Convergence (3)

$$\min_{x \in \mathbb{R}^2} f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

- **Modified Newton** (Armijo line-search;
if Hessian is < 0 switch to steepest descent)



Check Yourself

- What does the term complexity analysis refer to?
- What is the difference of analytical and arithmetic complexity
- Which method has better analytical complexity: Newton vs. steepest descent?

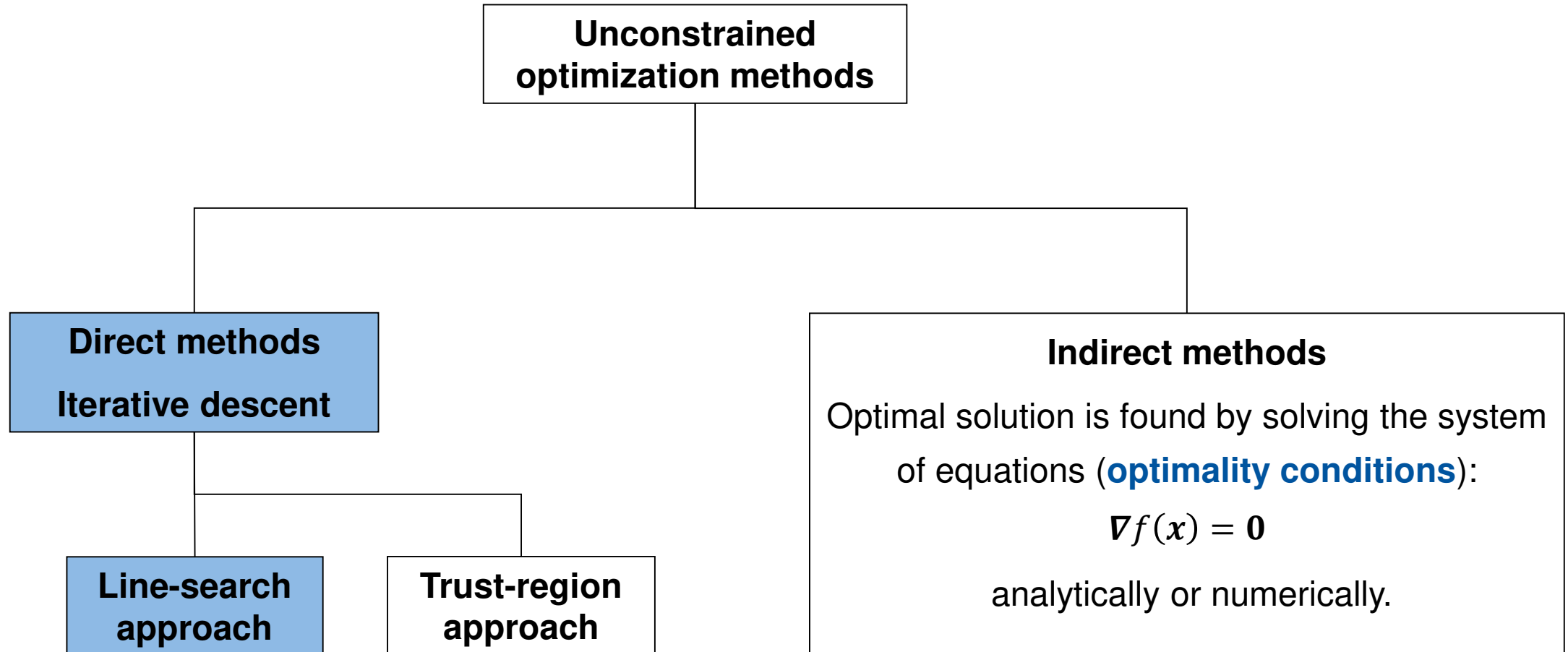


Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

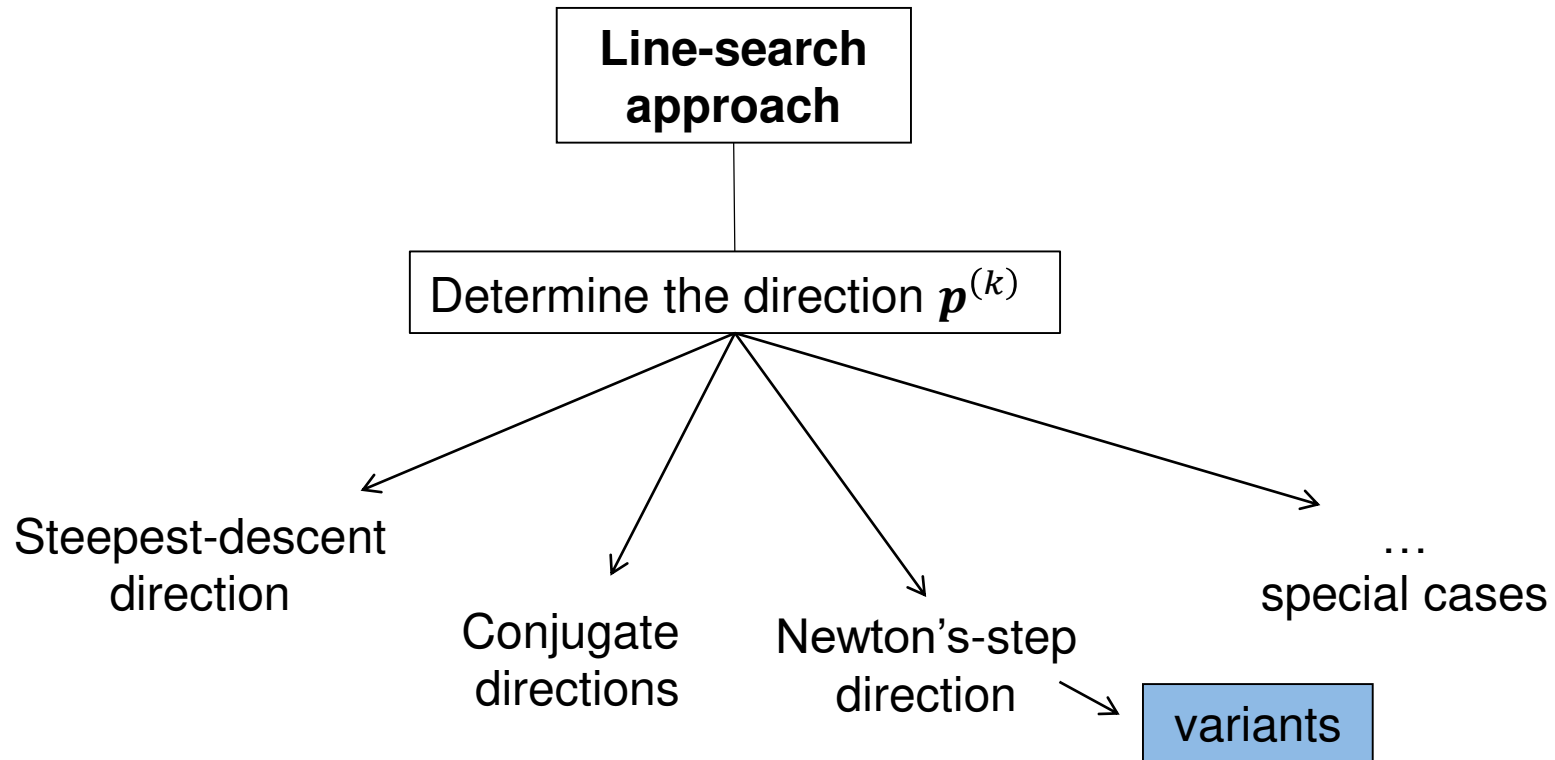
Line search: advanced directions

Solution Methods for Unconstrained Optimization



Determination of a Descent Direction: A Toolbox

Line-search approaches differ from each other with respect to the determination of descent direction and step length.



Inexact Newton Method (1)

Define: $f^{(k)} = f(x^{(k)})$ and $g^{(k)} := \nabla f(x^{(k)})$ and $H^{(k)} := \nabla^2 f(x^{(k)})$

From Newton's method:

$$x^{(k+1)} = x^{(k)} + p^{(k)}$$

$$[\nabla^2 f(x^{(k)})]p^{(k)} = -\nabla f(x^{(k)}) \Rightarrow H^{(k)}p^{(k)} = -g^{(k)}$$

Idea:

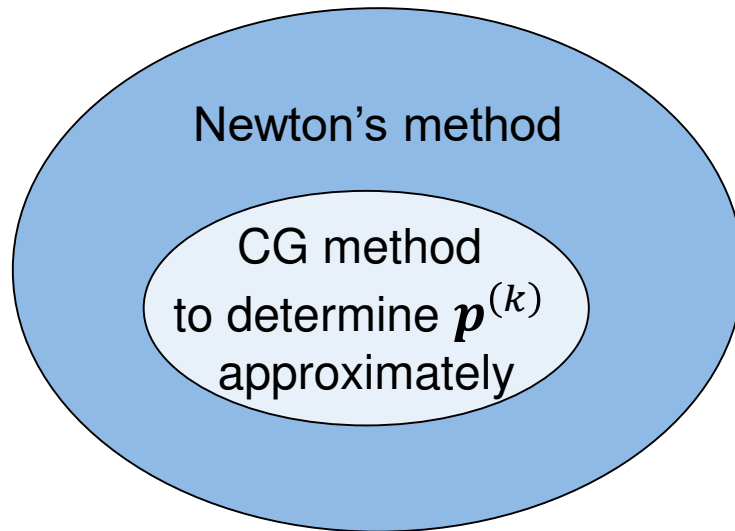
- The linear equation system, $H^{(k)}p^{(k)} = -g^{(k)}$, is solved approximately by an **iterative method**, e.g., by CG (conjugate gradients) **if $H^{(k)}$ is positive definite**.

Comments:

- LU- or Cholesky-decomposition – very high computational effort!
- Large errors occur for ill-conditioned problems.
- The **exact solution is not needed**.

Inexact Newton Method (2)

Newton-CG method:



Algorithm:

choose $\mathbf{x}^{(0)}$

for $k=0, 1, \dots$

if $\|\nabla f(\mathbf{x}^{(k)})\| \leq \varepsilon$ stop, else

calculate $\mathbf{g}^{(k)} := \nabla f(\mathbf{x}^{(k)})$ and $\mathbf{H}^{(k)} := \nabla^2 f(\mathbf{x}^{(k)})$

solve $\mathbf{H}^{(k)} \mathbf{p}^{(k)} = -\mathbf{g}^{(k)}$ for $\mathbf{p}^{(k)}$ with CG method

set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$

end for

Some line search strategy is needed. (Why?)

Modified Newton Method

Motivation: What if

- $\mathbf{H}^{(k)}$ is singular or almost singular (poorly conditioned)?
- $\mathbf{H}^{(k)}$ is not positive definite?

$$\mathbf{H}^{(k)} := \nabla^2 f(\mathbf{x}^{(k)})$$

| Idea | Approximations |
|--|---|
| <p>replace $\mathbf{H}^{(k)}$ by the approximation $\mathbf{B}^{(k)} \approx \mathbf{H}^{(k)}$</p> <p>→ $\mathbf{B}^{(k)} \mathbf{p}^{(k)} = -\mathbf{g}^{(k)}$</p> <p>$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$, ($\alpha_k$ from the line-search)</p> | <p>$\mathbf{B}^{(k)} = \mathbf{H}^{(k)} + \mathbf{E}^{(k)}$ with $\mathbf{E}^{(k)} = \tau_k \mathbf{I}$, $\tau_k \geq 0$ smartly chosen converges to steepest descent for $\tau_k \rightarrow \infty$</p> |

Alternatives exist, e.g., see [1]

Quasi-Newton Methods (1)

Idea: Reduce complexity by simplified calculation of $\mathbf{H}^{(k)}$ (Davidon):

- replace $\mathbf{H}^{(k)}$ by an approximation $\mathbf{B}^{(k)}$.
- instead of calculating $\mathbf{B}^{(k)}$, we look for a simple update using information from the last iterations.

$$\mathbf{H}^{(k)} := \nabla^2 f(\mathbf{x}^{(k)})$$

$$\mathbf{g}^{(k)} := \nabla f(\mathbf{x}^{(k)})$$

$$f^{(k)} := f(\mathbf{x}^{(k)})$$

Approach:

- Consider quadratic approximation of f at $\mathbf{x}^{(k)}$, $m^{(k)}(\mathbf{p}) = f^{(k)} + \mathbf{g}^{(k)T} \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{B}^{(k)} \mathbf{p}$.
- First order optimality condition: $\mathbf{p}^{(k)} = -\mathbf{B}^{(k)-1} \mathbf{g}^{(k)}$
- By convexity necessary and sufficient for minimization of $m^{(k)}(\mathbf{p})$.

symmetric positive definite

- Construct the quadratic approximation at $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$,

$$m^{(k+1)}(\mathbf{p}) = f^{(k+1)} + \mathbf{g}^{(k+1)T} \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{B}^{(k+1)} \mathbf{p}$$

- What **conditions** must $\mathbf{B}^{(k+1)}$ satisfy?

Quasi-Newton Methods (2)

Conditions on $\mathbf{B}^{(k+1)}$:

1. Gradient of $m^{(k+1)}$ at $\mathbf{x}^{(k)}$ and $\mathbf{x}^{(k+1)}$ must be equal to gradient of f .

| $\nabla m^{(k+1)}(\mathbf{p}) = \mathbf{g}^{(k+1)} + \mathbf{B}^{(k+1)}\mathbf{p}$ | |
|---|--|
| At $\mathbf{x} = \mathbf{x}^{(k+1)}$, $\mathbf{p} = 0$ We want $\nabla m^{(k+1)}(0) = \mathbf{g}^{(k+1)}$ | At $\mathbf{x} = \mathbf{x}^{(k)}$, $\mathbf{p} = -\alpha_k \mathbf{p}^{(k)}$ We want $\nabla m^{(k+1)}(-\alpha_k \mathbf{p}^{(k)}) = \mathbf{g}^{(k)}$ |
| Automatically satisfied | $\Rightarrow \mathbf{g}^{(k+1)} - \alpha_k \mathbf{B}^{(k+1)} \mathbf{p}^{(k)} = \mathbf{g}^{(k)}$ $\Rightarrow \mathbf{B}^{(k+1)} \underbrace{\alpha_k \mathbf{p}^{(k)}}_{=\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}} = \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}$ $\Rightarrow \boxed{\mathbf{B}^{(k+1)} \mathbf{s}^{(k)} = \mathbf{y}^{(k)}}$, where $\mathbf{s}^{(k)} = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$ and $\mathbf{y}^{(k)} = \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}$ |

2. Since, $\mathbf{B}^{(k+1)}$ is symmetric positive definite: $\mathbf{s}^{(k)T} \mathbf{B}^{(k+1)} \mathbf{s}^{(k)} > 0, \forall \mathbf{s}^{(k)} \neq 0 \Rightarrow \boxed{\mathbf{s}^{(k)T} \mathbf{y}^{(k)} > 0}$

★ Wolfe conditions (line-search) guarantee these constraints for all f , even when f is non-convex.

Quasi-Newton Methods (3)

Conditions on $\mathbf{B}^{(k+1)}$:

$\mathbf{B}^{(k+1)}\mathbf{s}^{(k)} = \mathbf{y}^{(k)}$ gives many solutions for $\mathbf{B}^{(k+1)}$

- **Unique** solution: $\mathbf{B}^{(k+1)}$ should be close to $\mathbf{B}^{(k)}$

$\min_{\mathbf{B}} \|\mathbf{B} - \mathbf{B}^{(k)}\|_W \leftarrow \text{weighted Frobenius-Norm}$

s. t. $\mathbf{B}^T = \mathbf{B}$

$\mathbf{B}\mathbf{s}^{(k)} = \mathbf{y}^{(k)}$

$\|A\|_W = \|W^{1/2}AW^{1/2}\|_F$, for any W s.t. $Wy_k = s_k$

$\|\cdot\|_F^2: R^{n \times n} \rightarrow R_{\geq 0}, \|C\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n c_{ij}^2$

$$\Rightarrow \mathbf{B}^{(k+1)} = \left(\mathbf{I} - \frac{1}{\mathbf{y}^{(k)T} \mathbf{s}^{(k)}} \mathbf{y}^{(k)} \mathbf{s}^{(k)T} \right) \mathbf{B}^{(k)} \left(\mathbf{I} - \frac{1}{\mathbf{y}^{(k)T} \mathbf{s}^{(k)}} \mathbf{s}^{(k)} \mathbf{y}^{(k)T} \right) + \frac{1}{\mathbf{y}^{(k)T} \mathbf{s}^{(k)}} \mathbf{y}^{(k)} \mathbf{y}^{(k)T} \rightarrow \text{DFP formula}$$

$$\Rightarrow \mathbf{B}^{(k+1)^{-1}} = \left(\mathbf{I} - \frac{1}{\mathbf{y}^{(k)T} \mathbf{s}^{(k)}} \mathbf{s}^{(k)} \mathbf{y}^{(k)T} \right) \mathbf{B}^{(k)^{-1}} \left(\mathbf{I} - \frac{1}{\mathbf{y}^{(k)T} \mathbf{s}^{(k)}} \mathbf{y}^{(k)} \mathbf{s}^{(k)T} \right) + \frac{1}{\mathbf{y}^{(k)T} \mathbf{s}^{(k)}} \mathbf{s}^{(k)} \mathbf{s}^{(k)T} \rightarrow \text{BFGS formula}$$

Check Yourself

- Explain the inexact Newton method.
- What is the main idea of the modified and quasi-Newton methods? Why are these methods advantageous?
- Why is it necessary to introduce a step-length control mechanism (line-search) into modified and quasi-Newton methods?



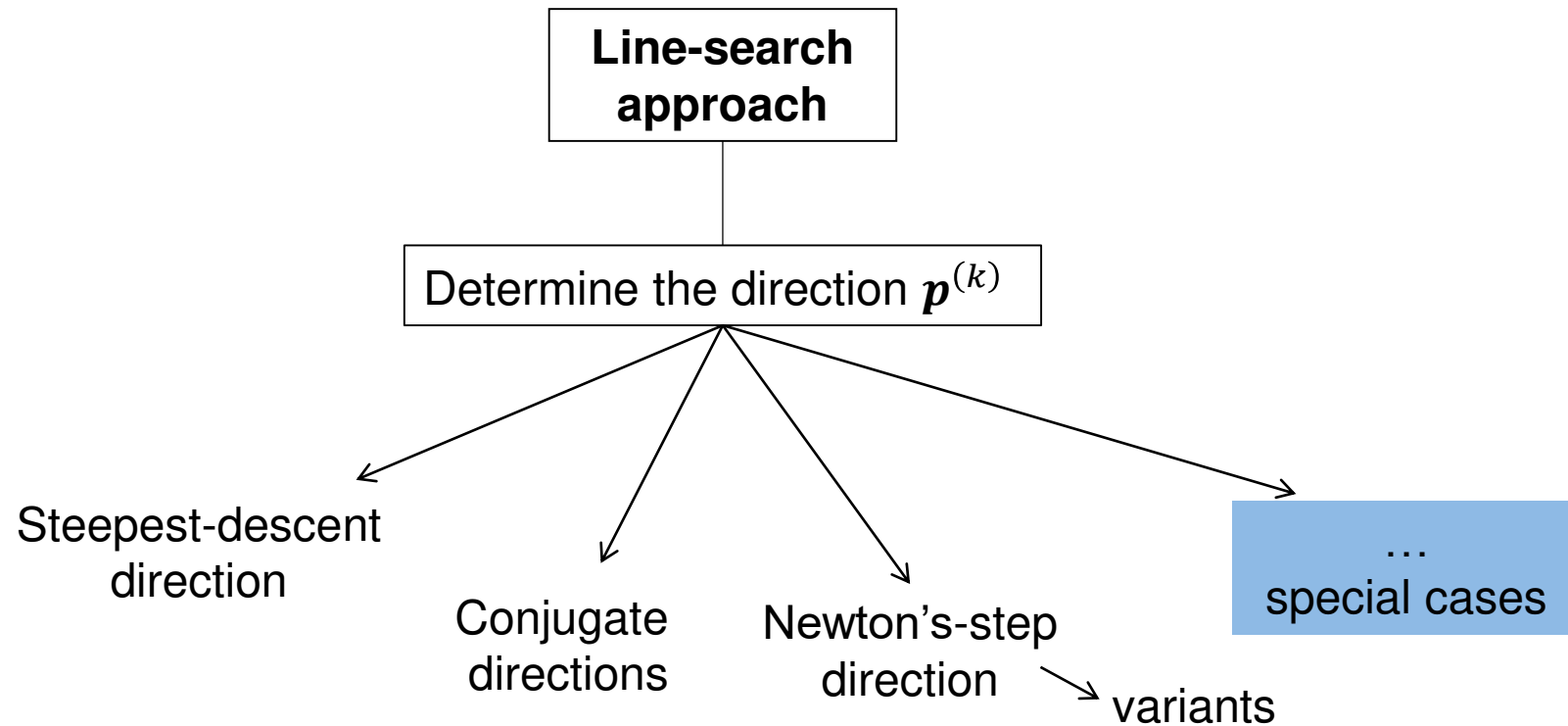
Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

Parameter estimation

Determination of a Descent Direction: A Toolbox

Line-search approaches differ from each other with respect to the determination of descent direction and step length.



Regression Problems: Least-Squares Formulation

Example:

Consider a batch reactor with the reaction $A \rightarrow B$ at constant temperature T_R . The reagent concentration c_A is measured at time instants t_j .

The reaction is of first order, therefore we can write the analytic solution:

$$\left. \frac{dc_A}{dt} \right|_t = -k \cdot c_A(t) \rightarrow c_A(t) = c_A|_{t=0} \cdot e^{-kt}$$

The reaction constant k and the reagent concentration at initial time $c_A(t=0)$ are unknown and should be determined from the measurements.

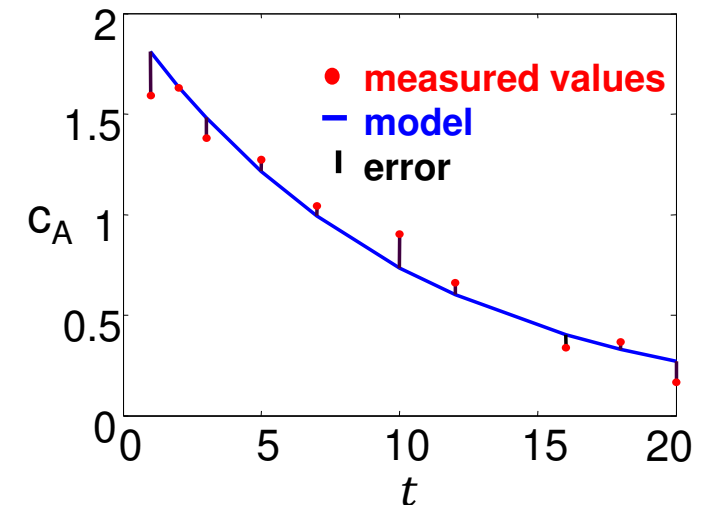
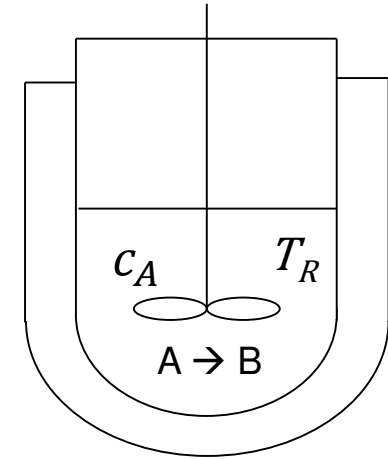
Optimization formulation uses $x_1 = c_A(t=0)$, $x_2 = k$:

$$c_{A, \text{theoretical}}(t_j) = \varphi(\mathbf{x}, t_j) = x_1 \cdot e^{x_2 \cdot t_j} \quad \text{model}$$

$$c_{A, \text{measured}}(t_j) = y_j \quad \text{measurement}$$

$$\varepsilon_j = y_j - \varphi(\mathbf{x}, t_j), \forall j = 1, \dots, m \quad \text{residual (error)}$$

$$\min_{\mathbf{x} \in \mathbb{R}^2} \frac{1}{2} \|\boldsymbol{\varepsilon}(\mathbf{x})\|_2^2 = \frac{1}{2} \sum_j \left(y_j - \phi(\mathbf{x}, t_j) \right)^2$$



Gauss-Newton Method

- $\min_{x \in \mathbb{R}^2} f(x) = \min_{x \in \mathbb{R}^2} \frac{1}{2} \|\varepsilon(x)\|_2^2 = \min_{x \in \mathbb{R}^2} \frac{1}{2} \varepsilon(x)^T \varepsilon(x)$

- Define: $J(x) := \nabla \varepsilon(x) \in \mathbb{R}^{m \times 2}$

$$\Rightarrow \nabla f(x) = J(x)^T \varepsilon(x)$$

$$\Rightarrow \nabla^2 f(x) = J(x)^T J(x) + \sum_{j=1}^m \varepsilon_j(x) \nabla^2 \varepsilon_j(x)$$

- The Hessian can be **approximated** by the first term in case of almost linear problems (i.e., $\nabla^2 \varepsilon_j(x) = 0$) or good starting values (i.e., small $\varepsilon_j(x)$)

- Newton's direction: $\nabla^2 f(x^{(k)}) p_k = -\nabla f(x^{(k)})$

- With Hessian approximation: $J^{(k)T} J^{(k)} p^{(k)} = -J^{(k)T} \varepsilon^{(k)}$

Remarks on Gauss-Newton Method

- If $J^{(k)}$ has full-rank, $\mathbf{p}^{(k)}$ is always a descent direction

$$\mathbf{p}^{(k)T} \cdot \nabla f(\mathbf{x}^{(k)}) = \mathbf{p}^{(k)T} \cdot J^{(k)T} \boldsymbol{\varepsilon}^{(k)} = -\mathbf{p}^{(k)T} \cdot J^{(k)T} \cdot J^{(k)} \cdot \mathbf{p}^{(k)} = -\|J^{(k)} \cdot \mathbf{p}^{(k)}\|_2^2 < 0,$$

The inequality is strict unless $J^{(k)} \cdot \mathbf{p}^{(k)} = 0 \Leftrightarrow J^{(k)T} \boldsymbol{\varepsilon}^{(k)} = \nabla f_k = 0$. \leftarrow Optimum

- In descent-direction $\mathbf{p}^{(k)}$, the step-length is determined as per the Wolfe-conditions
- For linear models the Jacobian J matrix is constant.
- The condition of minimization corresponds to the normal equations
- If $J^{(k)}$ is singular or almost singular, the descent direction $\mathbf{p}^{(k)}$ is, usually, not reliable. The method converges very poorly. Quasi-Newton methods are therefore more efficient.
- It is a local method

Check Yourself

- Where can least-squares problems be applied?
- When is a least-squares problem linear or nonlinear?
- What is the key idea of the Gauss-Newton method?

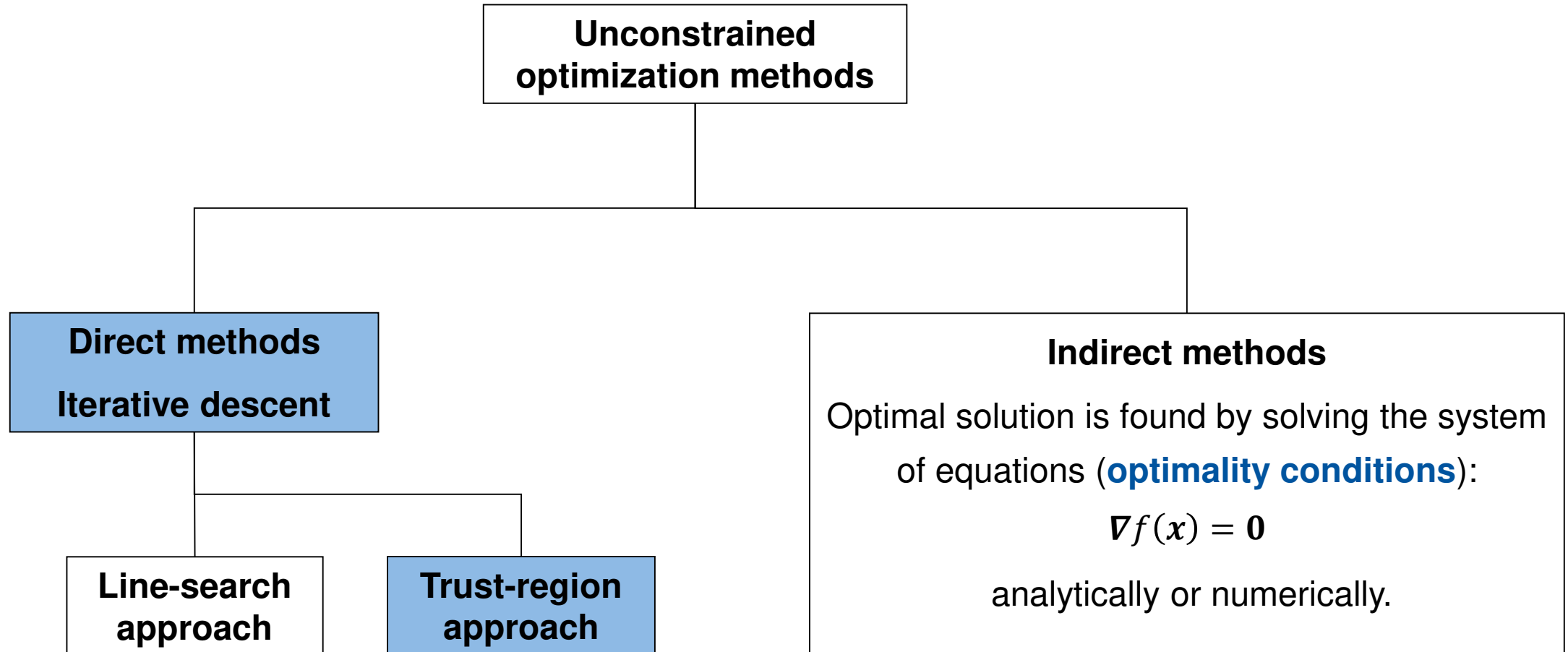


Applied Numerical Optimization

Prof. Alexander Mitsos, Ph.D.

Trust region method

Solution Methods for Unconstrained Optimization



Trust Region Method (1)

Idea:

- Approximate f at $\mathbf{x}^{(k)}$ by the quadratic model function $m^{(k)}$:

$$m^{(k)}(\mathbf{p}) = f^{(k)} + \mathbf{g}^{(k)T} \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{B}^{(k)} \mathbf{p}$$

where $f^{(k)} = f(\mathbf{x}^{(k)})$, $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ and $\mathbf{B}^{(k)}$ is symmetric

- Taylor-series: the approximation error is small for small \mathbf{p}
- For each iteration $k = 0, 1, \dots$ choose a **trust region radius** $\Delta^{(k)}$
- Solve the **minimization problem**: $\min_{\mathbf{p}} m^{(k)}(\mathbf{p})$ s. t. $\|\mathbf{p}\| \leq \Delta^{(k)}$ and set $\mathbf{p}^{(k)}$ to the solution found
- Set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^{(k)}$

Trust Region Method (2)

How to update the radius $\Delta^{(k)}$?

- Compare the agreement between the model function m_k and the objective function f at the previous iterations. Define **contraction rate** ρ_k as:

$$\rho_k = \frac{f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k)} + \mathbf{p}^{(k)})}{m^{(k)}(\mathbf{0}) - m^{(k)}(\mathbf{p}^{(k)})} = \frac{\text{actual reduction}}{\text{predicted reduction}}$$

As m_k is minimized over a domain containing $\mathbf{0}$:

$$m^{(k)}(\mathbf{0}) - m^{(k)}(\mathbf{p}^{(k)}) > 0$$

If $\rho_k < 0$! **reject** this step (ascent)

If $\rho_k \approx 1$! **increase** the radius: good agreement

If $\rho_k \approx 0$! **decrease** the radius: poor agreement

Trust Region Method (3)

Basic Algorithm:

choose $\Delta^{(max)} > 0, \Delta^{(0)} \in (0, \Delta^{(max)})$ and $\eta \in [0, \frac{1}{4})$

for $k = 0, 1, \dots$

calculate direction $\mathbf{p}^{(k)}$, contraction rate ρ_k

if $\rho_k < \frac{1}{4}$, $\Delta^{(k+1)} = \frac{\|\mathbf{p}^{(k)}\|}{4}$

else if $\rho_k > \frac{3}{4}$ and $\|\mathbf{p}^{(k)}\| = \Delta^{(k)}$, $\Delta^{(k+1)} = \min(2\Delta^{(k)}, \Delta^{(max)})$

else $\Delta^{(k+1)} = \Delta^{(k)}$

if $\rho_k > \eta$, $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^{(k)}$

else $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$

Trust Region Method (3)

Remarks:

- $\Delta^{(k)}$ is increased only if $\|p^{(k)}\|$ reaches the boundary of the domain.
- Strategies for the efficient solution of the minimization problem for $p^{(k)}$:
 - The Cauchy point: minimum along the steepest descent direction $(-g^{(k)})$, slow
 - The *Dogleg method*: applicable when $B^{(k)}$ is positive definite, fast (superlinear)
 - Steihaug's approach for large sparse matrices

The Dogleg Method

Idea:

- For a **large** $\Delta^{(k)}$: Newton step, $\mathbf{p}^{(k)} = \mathbf{p}^B = -\mathbf{B}^{(k)^{-1}} \mathbf{g}^{(k)}$. Where \mathbf{p}^B is the unconstrained minimum of m_k , $\|\mathbf{p}^B\| \leq \Delta^{(k)}$.
- For a **small** $\Delta^{(k)}$: search the solution along the direction $-\mathbf{g}^{(k)}$
- For an **intermediate** $\Delta^{(k)}$: additionally calculate $\mathbf{p}^U = -\frac{\mathbf{g}^{(k)^T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)^T} \mathbf{B}^{(k)} \mathbf{g}^{(k)}} \mathbf{g}^{(k)}$

Where \mathbf{p}^U is the unconstrained minimum of $m^{(k)}$ in the steepest descent direction.

- Then calculate the search-direction using a **linear combination** of

$$\mathbf{p}^U \text{ and } \mathbf{p}^B: \mathbf{p}^{(k)}(\tau) = \begin{cases} \tau \mathbf{p}^U & 0 \leq \tau \leq 1 \\ \mathbf{p}^U + (\tau - 1)(\mathbf{p}^B - \mathbf{p}^U) & 1 \leq \tau \leq 2 \end{cases}$$

with $\|\mathbf{p}^{(k)}(\tau^*)\| = \Delta^{(k)}$.

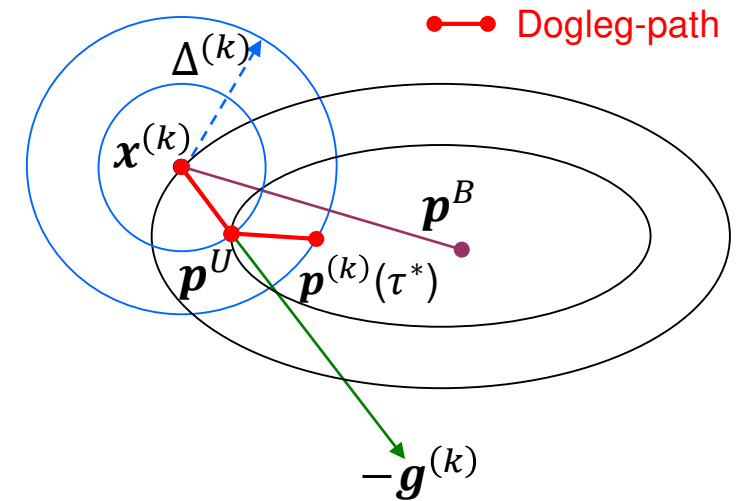


Illustration of Convergence – Rosenbrock Function

$$\min_{x \in \mathbb{R}^2} f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

- Solution point is $x^* = (1, 1)^T$ - why?

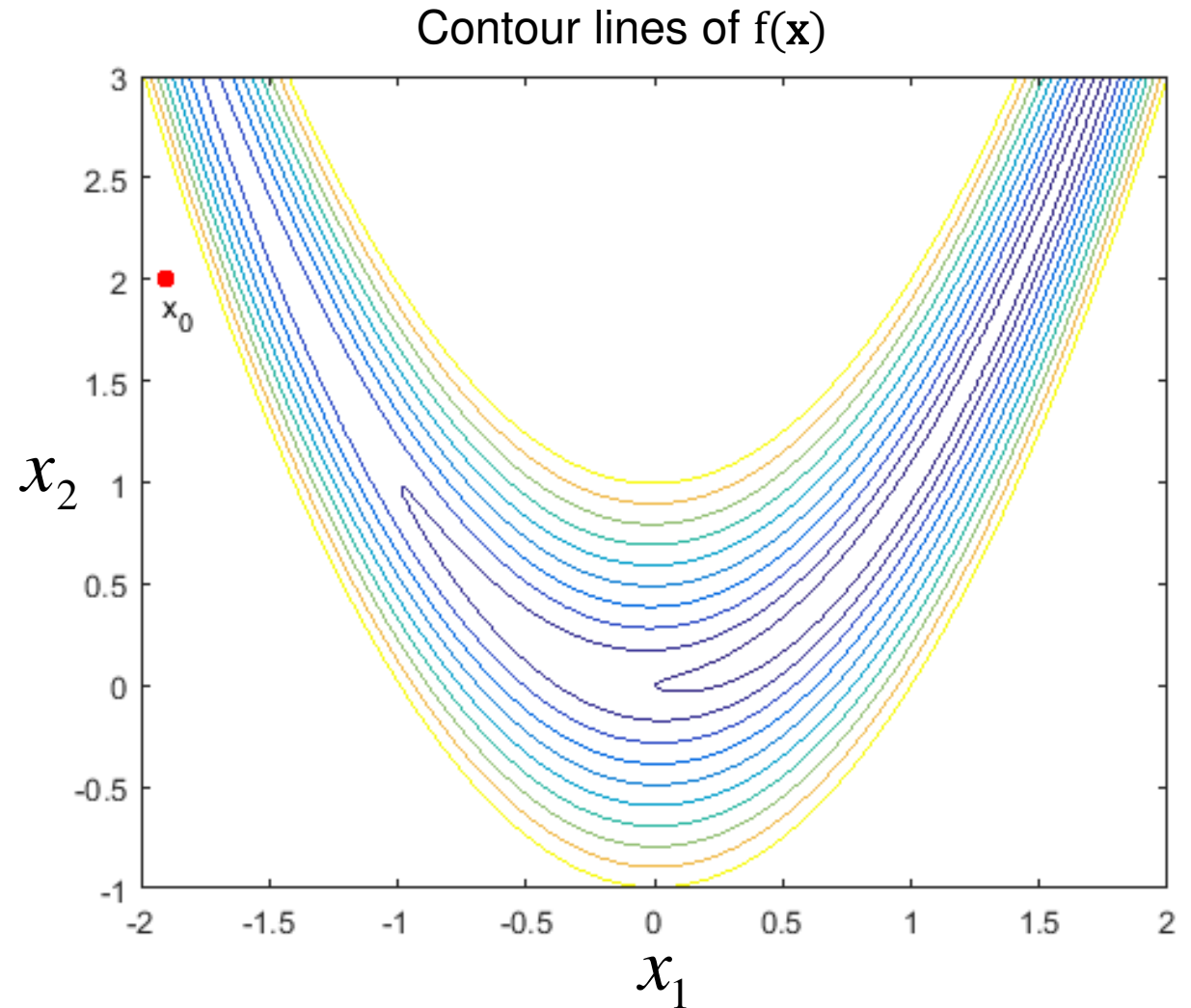
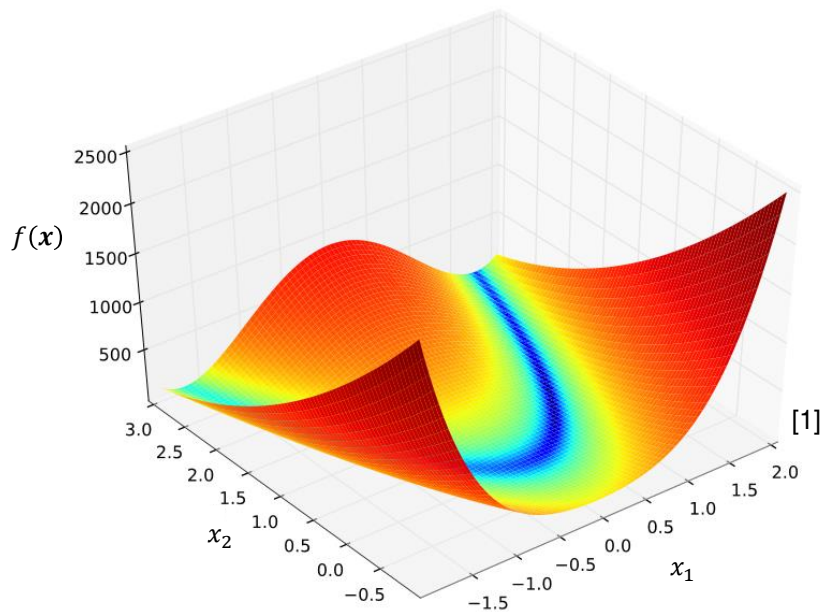
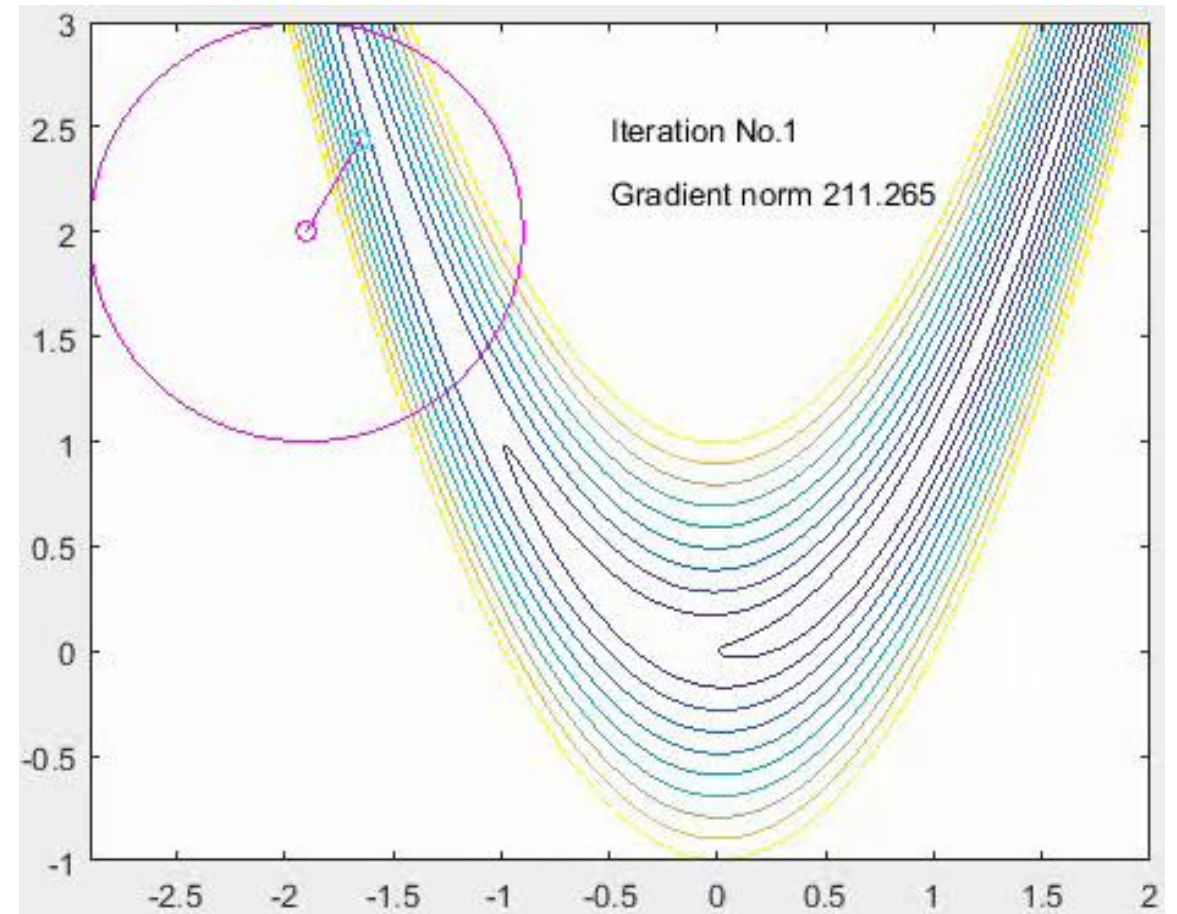
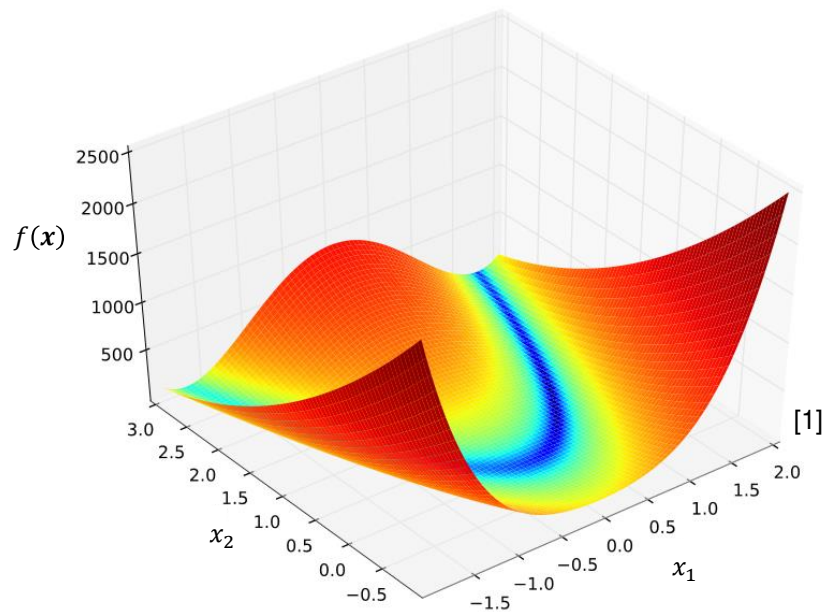


Illustration of Convergence – BFGS (Quasi-Newton Method)

$$\min_{x \in \mathbb{R}^2} f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

- Matlab trust-region (fminunc)



Check Yourself

- Explain trust region method.
- Which model problem is solved in trust-region methods?
- How is trust-region radius updated?