

# Introduction to Local and Global Optimization for NLP

Marco Trubian

Dipartimento di Scienze dell'Informazione (DSI)  
Università degli Studi di Milano



# Optimization

- Concerned with minimization/maximization of mathematical functions
- Often subject to constraints
- Important tool in analyzing physical, economic, chemical and biological systems
- Euler (1707-1783): **Nothing at all takes place in the universe in which some rule of the maximum or minimum does not apply.**
- Model  $\rightarrow$  apply algorithm  $\rightarrow$  check solution

# Unconstrained optimization

## Unconstrained

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in \mathbb{R}^n \end{array}$$

- $f$  : Objective function
- $\mathbf{x}$  : Decision variables
- $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- $f$  is assumed to be at least twice differentiable
- Examples:  $f(x) = 2x^3 - 3x^2$ .
- Important application: Data fitting and regression

# Constrained optimization

## Constrained

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, k \\ & h_i(\mathbf{x}) = 0, \quad i = 1, \dots, h \end{aligned}$$

- $g(\mathbf{x})$  : nonlinear inequality constraints
- $h(\mathbf{x})$  : nonlinear equality constraints
- Application: Resource constrained problems, transportation problems

# Convex optimization

## Convex

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in X \subseteq \mathbb{R}^n \end{array}$$

- $f$  :  $f$  is a convex function
- $X$  is convex: if  $\mathbf{x}^1, \mathbf{x}^2 \in X$  then  $(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2) \in X$ ,  $\lambda \in [0, 1]$
- Any local solution is global
- In constrained optimization,  $h(\mathbf{x}) = A\mathbf{x} - b$  and  $g(\mathbf{x})$  are convex
- Application: Controller design

# Discrete optimization

## Discrete

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in X \cap \mathbb{Z}^n \end{array}$$

- $\mathbf{x}$  can take only discrete values
- e.g.  $\mathbf{x} \in \{0, 1\}^n$
- Sometimes  $\mathbf{x}^1 \in \mathbb{R}, \mathbf{x}^2 \in \{0, 1\}^n$  implying mixed-integer model
- Application: facility location, routing, combinatorial problems

# Stochastic optimization

## Stochastic

$$\begin{aligned} \min \quad & Ef(\mathbf{x}; \omega) \\ \text{s.t.} \quad & \mathbf{x} \in \mathbb{R} \\ & c(\mathbf{x}; \omega) \geq 0, \omega \in \Omega \end{aligned}$$

- Random variable  $\omega$  belonging to a sample space  $\Omega$
- A constraint for each realization of uncertainty  $\omega$
- Minimize expected value
- Application: portfolio optimization problem

# Our focus

- Smooth nonlinearly constrained optimization problems
- Local solutions (satisfy optimality conditions)
- Main idea:
  - 1 Check if current point satisfies optimality conditions
  - 2 If not, obtain new iterate and return to 1.
- Finding new iterate requires using local information



# Curve fitting (1)

- Use a curve to fit experimental data
- Measurements  $y_1, \dots, y_m$  at times  $t_1, \dots, t_m$
- Model:  $\phi(t, \mathbf{x}) = x_1 + x_2 e^{-(x_3 - t)^2 / x_4} + x_5 \cos(x_6 t)$
- Model parameters:  $x_1, \dots, x_6$
- How does one set these parameters?

## Curve fitting (2)

- Define residuals  $r_j = y_j - \phi(t_j, \mathbf{x})$ ,  $j = 1, \dots, m$ ,
- Residuals: measure discrepancy between actual and estimated
- Specify performance metric:  $f(\mathbf{x}) = r_1^2(\mathbf{x}) + r_2^2(\mathbf{x}) + \dots + r_m^2(\mathbf{x})$

### Nonlinear Least Squares

$$\begin{aligned} \min \quad & \sum_{j=1}^m r_j^2 \\ \text{s.t.} \quad & \mathbf{x} \in \mathbb{R}^6 \end{aligned}$$

Weather forecasting: a limited number of parameters and a huge number of measurements

# Global and Local Minimizers

- Consider a function  $f(\mathbf{x})$  defined on  $X \subseteq \mathbb{R}^n$

**Definition** A point  $\mathbf{x}^* \in X$  is a **global minimizer** of  $f(\mathbf{x})$  if

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \forall \mathbf{x} \in X$$

$f(x) = (x - 2)^2$  has a global minimizer  $x^* = 2$

- Definition** A point  $\mathbf{x}^* \in X$  is a **local minimizer** of  $f(\mathbf{x})$  if there exists a neighborhood  $N(\mathbf{x}^*)$  of  $\mathbf{x}^*$  such that

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \forall \mathbf{x} \in N(\mathbf{x}^*)$$

- $f(x) = \sin(x)$  has local minimizers  $x^* = \{3\pi/2 \pm 2n\pi, n \in \mathbb{Z}\}$

# Strict and Isolated Minimizers

- Consider a function  $f(\mathbf{x})$  defined on  $X \subseteq \mathbb{R}^n$

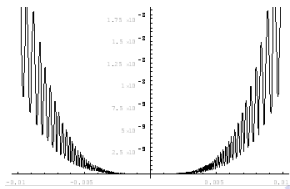
**Definition** A point  $\mathbf{x}^* \in X$  is a **strict local minimizer** of  $f(\mathbf{x})$  if there exists a neighborhood  $N(\mathbf{x}^*)$  of  $\mathbf{x}^*$  such that

$$f(\mathbf{x}^*) < f(\mathbf{x}) \quad \forall \mathbf{x} \in N(\mathbf{x}^*)$$

- $f(x) = \sin(x)$  has strict local minimizers  $x^* = \{3\pi/2 \pm 2n\pi, n \in \mathbb{Z}\}$

**Definition** A point  $\mathbf{x}^* \in X$  is a **isolated local minimizer** of  $f(\mathbf{x})$  if there exists a neighborhood  $N(\mathbf{x}^*)$  of  $\mathbf{x}^*$  such that  $\mathbf{x}^*$  is the only minimizer in  $N(\mathbf{x}^*)$

- $f(x) = x^4 \cos(1/x) + 2x^4$  has a strict global minimizer at  $x^* = 0$  but it is not isolated



# Our focus

- Smooth nonlinearly constrained optimization problems
- Local solutions (satisfy optimality conditions)
- Main idea:
  - ① Check if current point satisfies optimality conditions
  - ② If not, obtain new iterate and return to 1.
- Finding new iterate requires using local information

# The Fundamental Tools: Taylor's Theorems

**Theorem** (Taylor's First-Order Theorem) Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable and  $\mathbf{h} \in \mathbb{R}^n$ . Then we have that  $f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + o(\|\mathbf{h}\|^2)$ .

**Theorem** (Taylor's Second-Order Theorem) Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice continuously differentiable and  $\mathbf{h} \in \mathbb{R}^n$ . Then we have that  $f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T H(\mathbf{x}) \mathbf{h} + o(\|\mathbf{h}\|^3)$ .

In order to manipulate any possible kind of function  $f$  the algorithms locally approximate it with linear or quadratic models

$$l(\mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h}$$

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

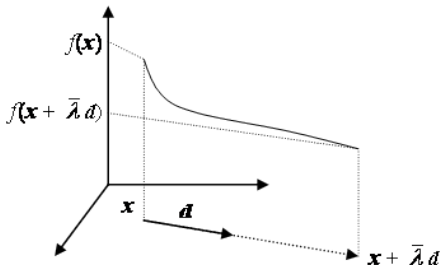
## Our focus

- Smooth nonlinearly constrained optimization problems
- Local solutions (satisfy optimality conditions)
- Main idea:
  - ① Check if current point satisfies optimality conditions
  - ② If not, obtain new iterate and return to 1.
- Finding new iterate requires using local information

# Descent directions: where to move

**Definition** (Descent direction and directional derivative) Given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  continuously differentiable in  $\mathbf{x}$  and a vector  $\mathbf{d}$

- if exists  $\bar{\lambda} > 0$  such that  $f(\mathbf{x} + \lambda \mathbf{d}) < f(\mathbf{x})$  for each  $0 < \lambda < \bar{\lambda}$  then  $\mathbf{d}$  is called a **descent direction** for  $f$  in  $\mathbf{x}$
- $\lim_{\lambda \rightarrow 0^+} \frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda} = \nabla f(\mathbf{x})^T \mathbf{d}$   
is called **directional derivative** of  $f$  in  $\mathbf{x}$  in the direction  $\mathbf{d}$





## Our focus

- Smooth nonlinearly constrained optimization problems
- Local solutions (satisfy optimality conditions)
- Main idea:
  - ① Check if current point satisfies optimality conditions
  - ② If not, obtain new iterate and return to 1.
- Finding new iterate requires using local information

# Optimality Conditions: when to stop

**Theorem** (First-Order Necessary Conditions) If  $\mathbf{x}^*$  is a local minimizer and  $f$  is continuously differentiable in a neighborhood of  $\mathbf{x}^*$ , then  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ .

**Theorem** (Second-Order Necessary Conditions) If  $\mathbf{x}^*$  is a local minimizer and the hessian matrix  $H(\mathbf{x})$  is continuous in a neighborhood of  $\mathbf{x}^*$ , then  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  and  $H(\mathbf{x}^*)$  is positive **semidefinite**.

**Theorem** (Second-Order Sufficient Conditions) Suppose that the hessian matrix  $H(\mathbf{x})$  is continuous in an neighborhood of  $\mathbf{x}^*$  and that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  and  $H(\mathbf{x}^*)$  is positive **definite**. Then  $\mathbf{x}^*$  is a **strict** local minimizer of  $f$ .

# Convex Programming (1)

**Definition** A square matrix  $H$  of order  $n$ , is called positive (semi)definite on a set  $X \subseteq \mathbb{R}^n$  if for each  $\mathbf{d} \in X$ ,  $\mathbf{d} \neq \mathbf{0}$ , then

$$\mathbf{d}^T H \mathbf{d} > 0 \quad \text{positive definite}$$

$$\mathbf{d}^T H \mathbf{d} \geq 0 \quad \text{positive semidefinite}$$

**Proposition** A **symmetric** matrix  $H$  is positive (semi)definite if and only if

- the determinants of all its principal minors are  $(\geq) > 0$
- its eigenvalues are  $(\geq) > 0$

**Proposition** Given a **symmetric** matrix  $H$  then

$$\lambda_{\min} \|\mathbf{d}\|^2 \leq \mathbf{d}^T H \mathbf{d} \leq \lambda_{\max} \|\mathbf{d}\|^2, \text{ for each } \mathbf{d},$$

where  $\lambda_{\min} := \min$  eigenvalue and  $\lambda_{\max} := \max$  eigenvalue of  $H$ .

**Observation** If  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  and the hessian matrix  $H(\mathbf{x}^*)$  is not positive (semi)definite then  $\mathbf{x}^*$  is neither a minimum nor a maximum.

# Convex Programming (2)

**Definition** A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  defined on  $X \subseteq \mathbb{R}^n$  is convex if  $X$  is convex and for each  $\mathbf{x}, \mathbf{y} \in X$  the following relationship holds

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \text{ for each } \lambda \in [0, 1].$$

**Proposition** Let  $X$  be a convex set and let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be twice continuously differentiable, then

- $f$  is convex if and only if the hessian matrix  $H(\mathbf{x})$  is positive semidefinite in  $\mathbf{x}$ , for each  $\mathbf{x} \in X$ .
- if the hessian matrix  $H(\mathbf{x})$  is positive **definite** in  $\mathbf{x}$  then  $f$  is **strictly** convex in a neighbourhood of  $\mathbf{x}$ .

## Theorem

- When  $f$  is a convex function, then any local minimizer  $\mathbf{x}^*$  is global.
- If  $f$  is a convex differentiable function, then any stationary point  $\mathbf{x}^*$  is a global minimizer of  $f$ .

# Quadratic Programming (1)

Optimization algorithms locally approximate  $f$  with quadratic models

$$\begin{aligned} \min \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x} \\ \text{t.c.} \quad & \mathbf{x} \in \mathbb{R}^n \end{aligned}$$

where  $Q$  is a symmetric square matrix of order  $n$ .

- $\nabla f(\mathbf{x}) = Q\mathbf{x} - \mathbf{b}$  (a linear function)
- $H(\mathbf{x}) = Q$  (a constant matrix)
- $Q\mathbf{x} = \mathbf{b}$  (first-order necessary conditions)

# Quadratic Programming (2)

Optimization algorithms locally approximate  $f$  with quadratic models

$$\begin{aligned} \min \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x} \\ \text{t.c.} \quad & \mathbf{x} \in \mathbb{R}^n \end{aligned}$$

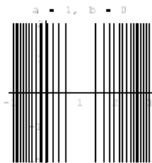
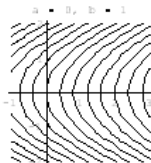
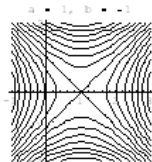
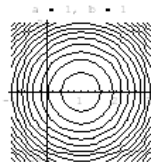
Possible cases:

- $Q$  is not positive semidefinite:  $f$  has not minima
- $Q$  is positive definite:  $\mathbf{x}^* = Q^{-1} \mathbf{b}$  is the only global minimizer
- $Q$  is positive semidefinite:
  - $Q$  is not singular:  $\mathbf{x}^* = Q^{-1} \mathbf{b}$  is the only global minimizer
  - $Q$  is singular:
    - no solutions or
    - infinite solutions

# Quadratic Programming (3)

Example:  $f(x, y) = \frac{1}{2} (ax^2 + by^2) - x$

$$f(\mathbf{x}) = \frac{1}{2} (x, y) \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} - (x, y) \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$



# Example of QP: Portfolio Optimization

- $r_i, i = 1, \dots, n$  r.v.: returns of  $n$  possible investments
- $\mu_i = E[r_i]$  and  $\sigma_i^2 = E[(r_i - \mu_i)^2]$  with normal distribution
- $\rho_{ij} = \frac{E[(r_i - \mu_i)(r_j - \mu_j)]}{\sigma_i \sigma_j}$  for  $i, j = 1, \dots, n$  correlations among investments pairs
- $x_i, i = 1, \dots, n$ , fraction of budget put into investment  $i$
- $\sum_{i=1}^n x_i = 1, \mathbf{x} \geq \mathbf{0}$ , all funds are invested
- $R = \sum_{i=1}^n x_i r_i$ , portfolio return
- $E[R] = \sum_{i=1}^n x_i E[r_i] = \mathbf{x}^T \boldsymbol{\mu}$ , expected return
- $Var[R] = E[(R - E[R])^2] = \sum_{i=1}^n \sum_{j=1}^n x_i x_j \sigma_i \sigma_j \rho_{ij} = \mathbf{x}^T G \mathbf{x}$   
where  $G_{ij} = \rho_{ij} \sigma_i \sigma_j$  is the symmetric positive semidefinite covariance matrix
- $\max \mathbf{x}^T \boldsymbol{\mu} - k \mathbf{x}^T G \mathbf{x}$ , s.t.  $\sum_{i=1}^n x_i = 1, \mathbf{x} \geq \mathbf{0}$ .



# Algorithm Overview (1)

Continuous optimization methods:

- 1 Given a starting point  $\mathbf{x}_0$ , generate an sequence  $\{\mathbf{x}_k\}_{k=0}^{\infty}$ .
- 2 Terminate the algorithm, when necessary conditions are satisfied with some accuracy, say  $\|\nabla f(\mathbf{x}_k)\| \leq \varepsilon$ .
- 3 Monotone algorithms requires that  $f(\mathbf{x}_k) < f(\mathbf{x}_{k-1})$  for all  $k$

How good is an optimization algorithm ?

# Local and Global Convergence

An optimization algorithm is *fair* if it **converges**...

**Definition** An algorithm is called **globally convergent** if it converges to a point  $\mathbf{x}^*$  s.t.  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  for each starting point  $\mathbf{x}_0 \in \mathbb{R}^n$

**Definition** An algorithm is called **locally convergent** if it converges to a point  $\mathbf{x}^*$  s.t.  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  only if the starting point  $\mathbf{x}_0 \in N(\mathbf{x}^*)$

**N.B.** nothing to do with the convergence to a local or a global optimum!

# Rate of Convergence

An optimization algorithm is **good** if it **converges rapidly!**

- Rate of convergence properties discuss the behavior of an algorithm close to a solution
- How fast does the algorithm converge?

Let  $\mathbf{x}_k$  be a sequence in  $\mathbb{R}^n$  that converges to  $\mathbf{x}^*$ . Convergence is:

- **Q-linear** if  $\exists r \in (0, 1)$  s.t.  $\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r$ , for  $k \geq \bar{k}$ .  
 E.g.:  $\{\mathbf{x}_k\} = \frac{1}{2^k}$
- **Q-superlinear** if  $\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0$ .  
 E.g.:  $\{\mathbf{x}_k\} = \frac{1}{k!}$
- **Q-quadratic** if  $\exists C > 0$  s.t.  $\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq C$ , for  $k \geq \bar{k}$ .  
 E.g.:  $\{\mathbf{x}_k\} = \frac{1}{2^{2^k}}$

Q-quadratically  $\rightarrow$  Q-superlinearly  $\rightarrow$  Q-linearly

## Algorithm Overview (2)

Continuous optimization methods:

- 1 Given a starting point  $\mathbf{x}_0$ , generate an sequence  $\{\mathbf{x}_k\}_{k=0}^{\infty}$ .
- 2 Terminate the algorithm, when necessary conditions are satisfied with some accuracy, say  $\|\nabla f(\mathbf{x}_k)\| \leq \varepsilon$ .
- 3 Monotone algorithms requires that  $f(\mathbf{x}_k) < f(\mathbf{x}_{k-1})$  for all  $k$

How does one determine  $\mathbf{x}_k$  given  $\mathbf{x}_{k-1}$ ?

# Two Main Strategies (1)

- **Line search:** Given an iterate  $\mathbf{x}_k$ , we first determine a direction  $\mathbf{h}_k$ . Then we determine a stepsize  $\alpha_k$  given by

$$\min_{\alpha_k} f(\mathbf{x}_k + \alpha_k \mathbf{h}_k)$$

The new iterate is given by  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{h}_k$ .

- **Trust region:** Construct a model function  $m_k$  using the information at  $\mathbf{x}_k$ . A trust-region radius  $\Delta_k$  is selected and we obtain a  $\mathbf{h}_k$  such that

$$\min_{\|\mathbf{h}_k\| \leq \Delta_k} m_k(\mathbf{h}_k)$$

If  $\mathbf{h}_k$  does not produce sufficient descent, we shrink  $\Delta_k$  and resolve.

## Two Strategies(2)

Line search:

- 1 Select direction
- 2 Determine stepsize

Trust Region:

- 1 Select trust region (stepsize)
- 2 Determine direction

# Line Search algorithms

## Descent Method;

```

{
  Choose  $\mathbf{x}_0 \in \mathbb{R}^n$ ;  $k := 0$ ;
  While  $\nabla f(\mathbf{x}_k) \neq \emptyset$ ;
  {
    compute  $\mathbf{d}_k \in \mathbb{R}^n$ ; /* descent direction */
    compute  $\alpha_k \in \mathbb{R}$ ; /* step along  $\mathbf{d}_k$  */
     $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{d}_k$ ;
     $k := k + 1$ ;
  }
}
    
```

# Line Search algorithms

## Descent Method;

```
{
  Choose  $\mathbf{x}_0 \in \mathbb{R}^n$ ;  $k := 0$ ;
  While  $\nabla f(\mathbf{x}_k) \neq \emptyset$ ;
  {
    compute  $\mathbf{d}_k \in \mathbb{R}^n$ ; /* descent direction */
    compute  $\alpha_k \in \mathbb{R}$ ; /* step along  $\mathbf{d}_k$  */
     $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{d}_k$ ;
     $k := k + 1$ ;
  }
}
```



# Stepsize specification

We assume that in  $\mathbf{x}$  a descent direction  $\mathbf{d}$  has been given

## Stepsize Problem

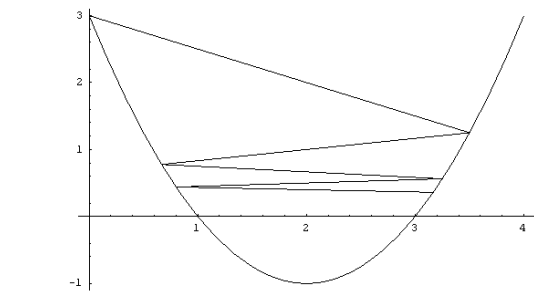
$$\min_{\alpha > 0} \phi(\alpha) = f(\mathbf{x} + \alpha \mathbf{d})$$

- Tradeoff between effort and accuracy
- Global minimizer would be too costly from a computational standpoint
- Exact linesearch:  $\alpha^*$  solves  $\phi'(\alpha) = \nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d} = 0$ .
- An inexact linesearch is generally used:
  - Cheap
  - Convergence rate does not rely on exact line search

# Observation

The simple decreasing of  $f$  is not enough.

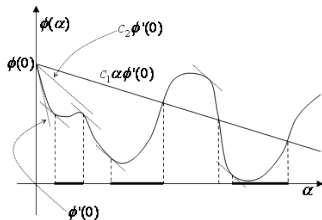
- $\phi(\alpha) = \alpha^2 - 4\alpha + 3$ , convex, with  $\alpha^* = 2$ .
- $\alpha_0 = 0$
- $\{\alpha_k\}$  generated by  $\alpha_k = 2 + (-1)^{k+1}(1 + 1/(k + 1))$
- $\phi(\alpha_{k+1}) < \phi(\alpha_k)$ , con  $k = 0, 2, \dots$



# Wolfe Conditions

To be effective inexact linesearch requires some conditions

- Sufficient Decrease:  $f(\mathbf{x} + \alpha \mathbf{d}) \leq f(\mathbf{x}) + c_1 \alpha \nabla f(\mathbf{x})^T \mathbf{d}$ ,  $c_1 \in (0, 1)$   
 $\phi(\alpha) \leq \phi(0) + \alpha c_1 \phi'(0)$
- Curvature condition:  $\nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d} \geq c_2 \nabla f(\mathbf{x})^T \mathbf{d}$ ,  $c_1 \in (0, 1)$   
 $\phi'(\alpha) \geq c_2 \phi'(0)$
- Collectively called **Wolfe conditions** where  $0 < c_1 < c_2 < 1$
- Strong Wolfe conditions introduce a sign constraint on curvature  
 $\|\nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d}\| \leq c_2 \|\nabla f(\mathbf{x})^T \mathbf{d}\|$



# Existence of stepsize

**Proposition** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable. Let  $\mathbf{d}$  be a descent direction at  $\mathbf{x}$  and assume that  $\phi(\alpha)$  is bounded below for  $\alpha > 0$ . Then if  $0 < c_1 < c_2 < 1$ , there exist steplengths satisfying the (strong) Wolfe conditions.

Looking for  $\alpha$ : backtracking line search

**Armijo's method;**

{

Choose  $\alpha_0 \in \mathbb{R}; \quad \alpha := \alpha_0;$

**While**  $f(\mathbf{x}_k + \alpha \mathbf{d}_k) > f(\mathbf{x}_k) + \alpha c_1 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k;$

$\alpha := \sigma \alpha; /* \text{backtracking} */$

$\alpha^* := \alpha;$

}

$\sigma \approx 0.9$ ,  $\alpha_0$  often set to 1 for Newton and Quasi-Newton methods

# Exact line search for QP

Let us consider a Quadratic Programming model

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

- $\phi(\alpha) = f(\mathbf{x} + \alpha \mathbf{d}) = \frac{1}{2} (\mathbf{x} + \alpha \mathbf{d})^T Q (\mathbf{x} + \alpha \mathbf{d}) - \mathbf{b}^T (\mathbf{x} + \alpha \mathbf{d})$
- $\phi'(\alpha) = \nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d} = 0.$
- $(Q(\mathbf{x} + \alpha \mathbf{d}) - \mathbf{b})^T \mathbf{d} = 0$
- $\alpha^* = -\frac{\mathbf{x}^T Q \mathbf{d} - \mathbf{b}^T \mathbf{d}}{\mathbf{d}^T Q \mathbf{d}} = -\frac{\nabla f(\mathbf{x})^T \mathbf{d}}{\mathbf{d}^T Q \mathbf{d}}$

This result is used in the convergence analysis of optimization algorithms.

# On computing $\alpha$

Other techniques for computing  $\alpha$

- Interpolation: using quadratic or cubic models
- Derivative free techniques, assuming convexity of  $\phi(\alpha)$ 
  - Golden Section Method,
  - Fibonacci Method,
  - Bisection Method.

# Convergence of inexact linesearch schemes (1)

If we define  $\theta_k$  as the angle between  $\mathbf{d}_k$  and  $-\nabla f_k$ , then

$$\cos \theta_k = \frac{-\nabla f(\mathbf{x})_k^T \mathbf{d}_k}{\|\nabla f(\mathbf{x})_k\| \cdot \|\mathbf{d}_k\|}$$

**Theorem** Let  $\mathbf{d}_k$  be a descent direction and let  $\alpha_k$  satisfy the Wolfe conditions. Also assume that  $f$  is bounded below on  $\mathbb{R}^n$  and continuously differentiable on  $N$  which contains the level set  $L_f := \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ , where  $\mathbf{x}_0$  is the starting point. We also assume that  $\nabla f$  is Lipschitz continuous on  $N$ . Then

$$\sum_{j=0}^{\infty} \cos^2 \theta_j \|\nabla f(\mathbf{x}_j)\|^2 < \infty.$$

# Convergence of inexact linesearch schemes (2)

**Proof:**

- $\nabla f_{k+1}^T \mathbf{d}_k \geq c_2 \nabla f_k^T \mathbf{d}_k$  (Curvature condition)
- $(\nabla f_{k+1} - \nabla f_k)^T \mathbf{d}_k \geq (c_2 - 1) \nabla f_k^T \mathbf{d}_k$
- $(\nabla f_{k+1} - \nabla f_k)^T \mathbf{d}_k \leq \|\nabla f_{k+1} - \nabla f_k\| \|\mathbf{d}_k\| \leq \alpha_k L \|\mathbf{x}_{k+1} - \mathbf{x}_k\| \|\mathbf{d}_k\| = \alpha_k L \|\mathbf{d}_k\|^2$  (from Schwartz and Lipschitz inequalities)
- $\alpha_k \geq \frac{c_2 - 1}{L} \frac{\nabla f_k^T \mathbf{d}_k}{\|\mathbf{d}_k\|^2}$
- $f_{k+1} \leq f_k + c_1 \alpha_k \mathbf{d}_k^T \nabla f_k = f_k - c_1 \frac{1 - c_2}{L} \frac{(\nabla f_k^T \mathbf{d}_k)^2}{\|\mathbf{d}_k\|^2}$  (Sufficient Decrease)
- $f_{k+1} \leq f_k - c \cos^2 \theta_k \|\nabla f_k\|^2$  where  $c = c_1(1 - c_2)/L$
- $f_{k+1} \leq f_0 - c \sum_{j=0}^k \cos^2 \theta_j \|\nabla f_j\|^2$  (by recursion)
- $c \sum_{j=0}^k \cos^2 \theta_j \|\nabla f_j\|^2 \leq f_0 - f_{k+1} < \infty$  (by boundness of  $f$ )  
 Zoutendijk condition



# Convergence of inexact linesearch schemes (3)

Zoutendijk condition

$$\sum_{j=0}^k \cos^2 \theta_j \|\nabla f_j\|^2 < \infty$$

implies

$$\cos^2 \theta_j \|\nabla f(\mathbf{x}_j)\|^2 \rightarrow 0$$

Hence if the algorithm satisfies **also** the **angle condition**

$$\cos \theta_k \geq \varepsilon > 0$$

then it **converges**

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\| = 0$$

## Convergence of inexact linesearch schemes (4)

Given a function  $f$  bounded below on  $\mathbb{R}^n$  and continuously differentiable on  $N$  which contains the level set  $L_f := \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ , with  $\nabla f$  Lipschitz continuous on  $N$  then an iterative method,

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

starting at  $\mathbf{x}_0$  **converges**, i.e.

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\| = 0$$

- if  $\mathbf{d}_k$  is a descent direction which satisfies the **angle condition**
- and  $\alpha_k$  satisfies the **Wolfe conditions**
  - Sufficient Decrease:  $f(\mathbf{x} + \alpha \mathbf{d}) \leq f(\mathbf{x}) + c_1 \alpha \nabla f(\mathbf{x})^T \mathbf{d}$ ,  $c_1 \in (0, 1)$
  - Curvature condition:  $\nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d} \geq c_2 \nabla f(\mathbf{x})^T \mathbf{d}$ ,  $c_1 \in (0, 1)$   
 where  $0 < c_1 < c_2 < 1$