# Unconstrained minimization
## Lectures for PHD course on
## Numerical Optimization

Enrico Bertolazzi

DII – Università di Trento

Notes

# Outline

Notes

Given $f : \mathbb{R}^n \mapsto \mathbb{R}$:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x)$$

the following regularity about $f(x)$ is assumed in the following:

## Assumption (Regularity assumption)

We assume $f \in C^1(\mathbb{R}^n)$ with Lipschitz continuous gradient, i.e. there exists $\gamma > 0$ such that

$$\left\| \nabla f(x)^T - \nabla f(y)^T \right\| \leq \gamma \left\| x - y \right\|, \qquad \forall x, y \in \mathbb{R}^n$$

Notes

> **Definition (Global minimum)**
>
> Given $f : \mathbb{R}^n \mapsto \mathbb{R}$ a point $\boldsymbol{x}_\star \in \mathbb{R}^n$ is a *global minimum* if
>
> $$f(\boldsymbol{x}_\star) \leq f(\boldsymbol{x}), \qquad \forall \boldsymbol{x} \in \mathbb{R}^n.$$

> **Definition (Local minimum)**
>
> Given $f : \mathbb{R}^n \mapsto \mathbb{R}$ a point $\boldsymbol{x}_\star \in \mathbb{R}^n$ is a *local minimum* if
>
> $$f(\boldsymbol{x}_\star) \leq f(\boldsymbol{x}), \qquad \forall \boldsymbol{x} \in B(\boldsymbol{x}_\star; \delta).$$

Obviously a global minimum is a local minimum. Find a global minimum in general is not an easy task. The algorithms presented in the sequel will approximate local minima's.

Notes

## Definition (Strict global minimum)

*Given* $f : \mathbb{R}^n \mapsto \mathbb{R}$ *a point* $x_\star \in \mathbb{R}^n$ *is a* **strict global minimum** *if*

$$f(x_\star) < f(x), \qquad \forall x \in \mathbb{R}^n \setminus \{x_\star\}.$$

## Definition (Strict local minimum)

*Given* $f : \mathbb{R}^n \mapsto \mathbb{R}$ *a point* $x_\star \in \mathbb{R}^n$ *is a* **strict local minimum** *if*

$$f(x_\star) < f(x), \qquad \forall x \in B(x_\star; \delta) \setminus \{x_\star\}.$$

Obviously a strict global minimum is a strict local minimum.

Notes

# First order Necessary condition

## Lemma (First order Necessary condition for local minimum)

*Given* $f : \mathbb{R}^n \mapsto \mathbb{R}$ *satisfying the regularity assumption. If a point* $x_\star \in \mathbb{R}^n$ *is a local minimum then*

$$\nabla f(x_\star)^T = 0.$$

**Proof:** Consider a generic direction $d$, then for $\delta$ small enough we have

$$\lambda^{-1}\big(f(x_\star + \lambda d) - f(x_\star)\big) \leq 0, \qquad 0 < \lambda < \delta$$

so that

$$\lim_{\lambda \to 0} \lambda^{-1}\big(f(x_\star + \lambda d) - f(x_\star)\big) = \nabla f(x_\star)d \leq 0,$$

because $d$ is a generic direction we have $\nabla f(x_\star)^T = 0$. $\quad \square$

Notes

## Remark

1. *The first order necessary condition do not discriminate maximum, minimum, or saddle points.*

2. *To discriminate maximum and minimum we need more information, e.g. second order derivative of $f(x)$.*

3. *With second order derivative we can build necessary and sufficient condition for a minima.*

4. *In general using only first and second order derivative at the point $x_\star$ it is not possible to deduce a necessary and sufficient condition for a minima.*

Notes

# Second order Necessary condition

**Lemma (Second order Necessary condition for local minimum)**

*Given* $f \in C^2(\mathbb{R}^n)$ *if a point* $\boldsymbol{x}_\star \in \mathbb{R}^n$ *is a local minimum then* $\nabla f(\boldsymbol{x}_\star)^T = \boldsymbol{0}$ *and* $\nabla^2 f(\boldsymbol{x}_\star)$ *is semi-definite positive, i.e.*

$$\boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}_\star) \boldsymbol{d} \geq 0, \qquad \forall \boldsymbol{d} \in \mathbb{R}^n$$

**Example**

This condition is only, necessary, in fact consider $f(\boldsymbol{x}) = x_1^2 - x_2^3$,

$$\nabla f(\boldsymbol{x}) = \left(2x_1, -3x_2^2\right), \quad \nabla^2 f(\boldsymbol{x}) = \begin{pmatrix} 2 & 0 \\ 0 & -6x_2 \end{pmatrix}$$

for the point $\boldsymbol{x}_\star = \boldsymbol{0}$ we have $\nabla f(\boldsymbol{0}) = \boldsymbol{0}$ and $\nabla^2 f(\boldsymbol{0})$ semi-definite positive, but $\boldsymbol{0}$ is a saddle point not a minimum.

Notes

**Proof:** The condition $\nabla f(\boldsymbol{x}_\star)^T = \boldsymbol{0}$ comes from first order necessary conditions. Consider now a generic direction $\boldsymbol{d}$, and the finite difference:

$$\frac{f(\boldsymbol{x}_\star + \lambda\boldsymbol{d}) - 2f(\boldsymbol{x}_\star) + f(\boldsymbol{x}_\star - \lambda\boldsymbol{d})}{\lambda^2} \geq 0$$

by using Taylor expansion for $f(\boldsymbol{x})$

$$f(\boldsymbol{x}_\star \pm \lambda\boldsymbol{d}) = f(\boldsymbol{x}_\star) \pm \nabla f(\boldsymbol{x}_\star)\lambda\boldsymbol{d} + \frac{\lambda^2}{2}\boldsymbol{d}^T\nabla^2 f(\boldsymbol{x}_\star)\boldsymbol{d} + o(\lambda^2)$$

and from the previous inequality

$$\boldsymbol{d}^T\nabla^2 f(\boldsymbol{x}_\star)\boldsymbol{d} + 2o(\lambda^2)/\lambda^2 \geq 0$$

taking the limit $\lambda \to 0$ and form the arbitrariness of $\boldsymbol{d}$ we have that $\nabla^2 f(\boldsymbol{x}_\star)$ must be semi-definite positive. $\square$

Notes

# Second order sufficient condition

**Lemma (Second order sufficient condition for local minimum)**

*Given* $f \in C^2(\mathbb{R}^n)$ *if a point* $\boldsymbol{x}_\star \in \mathbb{R}^n$ *satisfy:*

1. $\nabla f(\boldsymbol{x}_\star)^T = \boldsymbol{0}$;
2. $\nabla^2 f(\boldsymbol{x}_\star)$ *is definite positive; i.e.*

$$\boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}_\star) \boldsymbol{d} > 0, \qquad \forall \boldsymbol{d} \in \mathbb{R}^n \setminus \{\boldsymbol{x}_\star\}$$

*then* $\boldsymbol{x}_\star \in \mathbb{R}^n$ *is a strict local minimum.*

**Remark**

*Because* $\nabla^2 f(\boldsymbol{x}_\star)$ *is symmetric we can write*

$$\lambda_{\min} \boldsymbol{d}^T \boldsymbol{d} \leq \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}_\star) \boldsymbol{d} \leq \lambda_{\max} \boldsymbol{d}^T \boldsymbol{d}$$

*If* $\nabla^2 f(\boldsymbol{x}_\star)$ *is positive definite we have* $\lambda_{\min} > 0$.

Notes

**Proof:** Consider now a generic direction $d$, and the Taylor expansion for $f(x)$

$$f(x_\star + d) = f(x_\star) + \nabla f(x_\star)d + \frac{1}{2}d^T \nabla^2 f(x_\star)d + o(\|d\|^2)$$

$$\geq f(x_\star) + \frac{1}{2}\lambda_{min}\|d\|^2 + o(\|d\|^2)$$

$$\geq f(x_\star) + \frac{1}{2}\lambda_{min}\|d\|^2 \left(1 + o(\|d\|^2)/\|d\|^2\right)$$

choosing $d$ small enough we can write

$$f(x_\star + d) \geq f(x_\star) + \frac{1}{4}\lambda_{min}\|d\|^2 > f(x_\star), \qquad d \neq 0, \ \|d\| \leq \delta.$$

i.e. $x_\star$ is a strict minimum.    □

Notes

Notes

# How to find a minimum

Given $f : \mathbb{R}^n \mapsto \mathbb{R}$:    $\text{minimize}_{x \in \mathbb{R}^n}$    $f(x)$.

1. We can solve the problem by solving the **necessary condition**. i.e by solving the nonlinear systems

$$\nabla f(x)^T = \mathbf{0}.$$

2. Using such an approach we looses the information about $f(x)$.

3. Moreover such an approach can find solution corresponding to a maximum or saddle points.

4. A better approach is to use all the information and try to build **minimizing procedure**, i.e. procedures that, starting from a point $x_0$ build a sequence $\{x_k\}$ such that $f(x_{k+1}) \leq f(x_k)$. In this way, at least, we avoid to converge to a **strict maximum**.

Notes

# Iterative Methods

- In practice, rarely we are able to provide an explicit minimizer.

- Iterative method: given starting **guess** $x_0$, generate the sequence,

$$\{x_k\}, \qquad k = 1, 2, \ldots$$

- **AIM:** ensure that (a subsequence) has some favorable limiting properties:
  - satisfies first-order necessary conditions
  - satisfies second-order necessary conditions

Notes

# Line-search Methods

A generic iterative minimization procedure can be sketched as follows:

- calculate a **search direction** $p_k$ from $x_k$
- ensure that this direction is a **descent direction**, i.e.

$$\nabla f(x_k) p_k < 0, \qquad \text{whenever } \nabla f(x_k)^T \neq 0$$

  so that, at least for small steps along $p_k$, the objective function $f(x)$ will be reduced

- use **line-search** to calculate a suitable step-length $\alpha_k > 0$ so that

$$f(x_k + \alpha_k p_k) < f(x_k).$$

- Update the point:

$$x_{k+1} = x_k + \alpha_k p_k$$

Notes

# Generic minimization algorithm

Written with a pseudo-code the minimization procedure is the following algorithm:

---

**Generic minimization algorithm**

Given an initial guess $x_0$, let $k = 0$;
**while** not converged **do**
　　Find a descent direction $p_k$ at $x_k$;
　　Compute a step size $\alpha_k$ using a line-search along $p_k$.
　　Set $x_{k+1} = x_k + \alpha_k p_k$ and increase $k$ by $1$.
**end while**

---

The crucial points which differentiate the algorithms are:

1. The computation of the direction $p_k$;

2. The computation of the step size $\alpha_k$.

Notes

# Practical Line-search methods

- The first developed minimization algorithms try to solve

$$\alpha_k = \arg\min_{\alpha>0} \mathsf{f}(\boldsymbol{x}_k + \alpha\boldsymbol{p}_k)$$

  - performing **exact line-search** by univariate minimization;
  - rather expensive and certainly not cost effective.

- Modern methods implements **inexact** line-search:
  - ensure steps are neither too long nor too short
  - try to pick **useful** initial step size for fast convergence
  - best methods are based on:
    - backtracking–Armijo search;
    - Armijo–Goldstein search;
    - Franke–Wolfe search;

Notes

# backtracking line-search

To obtain a monotone decreasing sequence we can use the following algorithm:

---

### Backtracking line-search

Given $\alpha_{\text{init}}$ (e.g., $\alpha_{\text{init}} = 1$);
Given $\tau \in (0, 1)$ typically $\tau = 0.5$;
Let $\alpha^{(0)} = \alpha_{\text{init}}$;
**while** not $f(\boldsymbol{x}_k + \alpha^{(\ell)}\boldsymbol{p}_k) < f(\boldsymbol{x}_k)$ **do**
    set $\alpha^{(\ell+1)} = \tau\alpha^{(\ell)}$;
    increase $\ell$ by $1$;
**end while**
Set $\alpha_k = \alpha^{(\ell)}$.

---

To be effective the previous algorithm should terminate in a finite number of steps. In the following we prove that if $\boldsymbol{p}_k$ is a descent direction then a slight modification of the algorithm will terminate.

Notes

# Existence of a descent step

## Lemma (Descent Lemma)

*Suppose that $f(\boldsymbol{x})$ satisfy the standard assumptions and that $\boldsymbol{p}_k$ is a descent direction at $\boldsymbol{x}_k$, i.e. $\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k < 0$. Then we have*

$$f(\boldsymbol{x}_k + \alpha\boldsymbol{p}_k) \leq f(\boldsymbol{x}_k) + \alpha\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k + \frac{\gamma}{2}\alpha^2 \|\boldsymbol{p}_k\|^2$$

*for all $\alpha \in [0, \alpha_k^\star]$ where* $\quad \alpha_k^\star = \dfrac{-2\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k}{\gamma \|\boldsymbol{p}_k\|^2} > 0$

## Assumption (Regularity assumption)

*We assume $f \in C^1(\mathbb{R}^n)$ with Lipschitz continuous gradient, i.e. there exists $\gamma > 0$ such that*

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\| \leq \gamma \|\boldsymbol{x} - \boldsymbol{y}\|, \qquad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n$$

Notes

**Proof:** Let be $g(\alpha) = \mathsf{f}(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k)$ then we can write:

$$g(\alpha) - g(0) = \int_0^\alpha g'(\xi)d\xi = \alpha g'(0) + \int_0^\alpha \left(g'(\xi) - g'(0)\right)d\xi$$

$$= \alpha \nabla \mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k + \int_0^\alpha \left(\nabla \mathsf{f}(\boldsymbol{x}_k + \xi \boldsymbol{p}_k) - \nabla \mathsf{f}(\boldsymbol{x}_k)\right)\boldsymbol{p}_k \, d\xi$$

$$\leq \alpha \nabla \mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k + \int_0^\alpha \|\nabla \mathsf{f}(\boldsymbol{x}_k + \xi \boldsymbol{p}_k) - \nabla \mathsf{f}(\boldsymbol{x}_k)\| \, \|\boldsymbol{p}_k\| \, d\xi$$

$$\leq \alpha \nabla \mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k + \|\boldsymbol{p}_k\|^2 \int_0^\alpha \gamma\xi \, d\xi$$

$$\leq \alpha \nabla \mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k + \frac{\gamma\alpha^2}{2}\|\boldsymbol{p}_k\|^2 = \alpha \left[\nabla \mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k + \frac{\gamma\alpha}{2}\|\boldsymbol{p}_k\|^2\right].$$
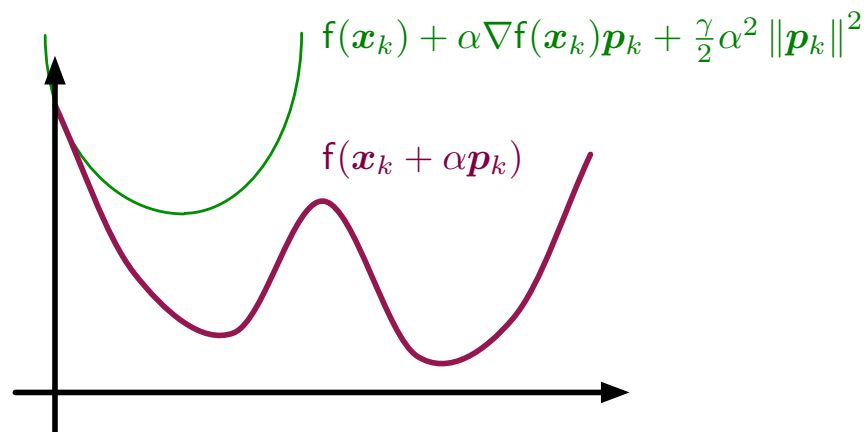
now the lemma follows trivially. $\square$

Notes

■ The **descent lemma** means that there is a parabola that is entirely over the function $f(x)$ in the direction $p_k$ if this is a descent direction.

■ The second part of the lemma permits to ensure a **minimal** reduction if the step length is chosen to be $\alpha_k = \alpha_k^\star/2$.

$$f(x_k) + \alpha \nabla f(x_k) p_k + \tfrac{\gamma}{2}\alpha^2 \left\| p_k \right\|^2$$

$$f(x_k + \alpha p_k)$$

Notes

# Descent direction failure

- The simple request to have a descent direction may be not enough.
- In fact, step length may be asymptotically too short
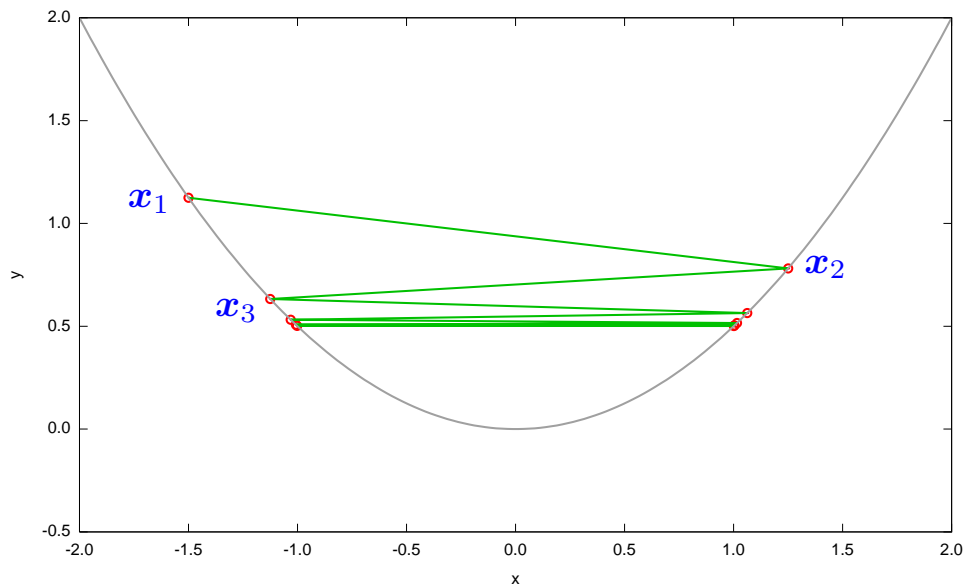- Or step length may be asymptotically too long

Notes

# Steps may be too long

The objective function is $f(x) = x^2$ and the iterates are generated by the descent directions $p_k = (-1)^{k+1}$ from $x_0 = 2$ with:

$$x_{k+1} = x_k + \alpha_k p_k, \qquad \alpha_k = 2 + 3 \cdot 2^{-(k+1)}$$
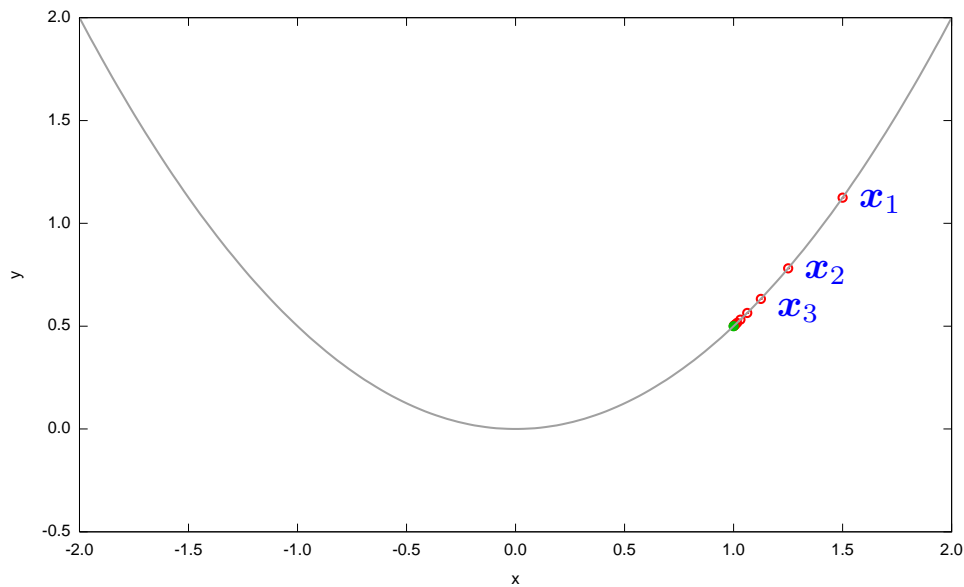
Notes

# Steps may be too short

The objective function is $f(x) = x^2$ and the iterates are generated by the descent directions $p_k = -1$ from $x_0 = 2$ with:

$$x_{k+1} = x_k + \alpha_k p_k, \qquad \alpha_k = 2^{-(k+1)}$$

Notes

Notes

# Armijo condition

To prevent large steps relative to the decreasing of $f(\boldsymbol{x})$ we require that

$$f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \leq f(\boldsymbol{x}_k) + \alpha_k \beta \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$

for some $\beta \in (0, 1)$. Typical values of $\beta$ ranges form $10^{-4}$ to $0.1$.

Notes

## Backtracking Armijo line-search

Given $\alpha_{\text{init}}$ (e.g., $\alpha_{\text{init}} = 1$);
Given $\tau \in (0, 1)$ typically $\tau = 0.5$;
Let $\alpha^{(0)} = \alpha_{\text{init}}$;
**while** not $f(\boldsymbol{x}_k + \alpha^{(\ell)} \boldsymbol{p}_k) \leq f(\boldsymbol{x}_k) + \alpha^{(\ell)} \beta \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$ **do**
    set $\alpha^{(\ell+1)} = \tau \alpha^{(\ell)}$;
    increase $\ell$ by $1$;
**end while**
Set $\alpha_k = \alpha^{(\ell)}$.

- Backtracking Armijo line-search prevents the step from getting too large.
- Now the question is: will the backtracking Armijo line-search terminate in a finite number of steps ?

Notes

# Finite termination of Armijo line-search

## Theorem (Finite termination of Armijo linesearch)

*Suppose that $f(\boldsymbol{x})$ satisfy the standard assumptions and $\beta \in (0,1)$ and that $\boldsymbol{p}_k$ is a descent direction at $\boldsymbol{x}_k$. Then the Armijo condition*

$$f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \leq f(\boldsymbol{x}_k) + \alpha_k \beta \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$

*is satisfied when $\alpha_k \in [0, \omega_k]$ where*   $\omega_k = \dfrac{2(\beta - 1)\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k}{\gamma \left\| \boldsymbol{p}_k \right\|^2}$

## Assumption (Regularity assumption)

*We assume $f \in C^1(\mathbb{R}^n)$ with Lipschitz continuous gradient, i.e. there exists $\gamma > 0$ such that*

$$\left\| \nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y}) \right\| \leq \gamma \left\| \boldsymbol{x} - \boldsymbol{y} \right\|, \qquad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n$$

Notes

# Finite termination of Armijo line-search

To prove finite termination we need the following Taylor expansion due to the regularity assumption:

$$\mathsf{f}(\boldsymbol{x} + \alpha\boldsymbol{p}) = \mathsf{f}(\boldsymbol{x}) + \alpha\nabla\mathsf{f}(\boldsymbol{x})\boldsymbol{p} + E \quad \text{where} \quad |E| \leq \frac{\gamma}{2}\alpha^2 \|\boldsymbol{p}\|^2$$

**Proof:** If $\alpha \leq \omega_k$ we have $\alpha\gamma \|\boldsymbol{p}_k\|^2 \leq 2(\beta - 1)\nabla\mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k$ and by using Taylor expansion

$$\mathsf{f}(\boldsymbol{x}_k + \alpha\boldsymbol{p}_k) \leq \mathsf{f}(\boldsymbol{x}_k) + \alpha\nabla\mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k + \frac{\gamma}{2}\alpha^2 \|\boldsymbol{p}_k\|^2$$

$$\leq \mathsf{f}(\boldsymbol{x}_k) + \alpha\nabla\mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k + \alpha(\beta - 1)\nabla\mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k$$

$$\leq \mathsf{f}(\boldsymbol{x}_k) + \alpha\beta\nabla\mathsf{f}(\boldsymbol{x}_k)\boldsymbol{p}_k$$

$\square$

Notes

# Finite termination of Armijo line-search

## Corollary (Finite termination of Armijo linesearch)

*Suppose that* $f(x)$ *satisfy the standard assumptions and* $\beta \in (0, 1)$ *and that* $p_k$ *is a descent direction at* $x_k$. *Then the step-size generated by then backtracking-Armijo line-search terminates with*

$$\alpha_k \geq \min\{\alpha_{\text{init}}, \tau\omega_k\}, \qquad \omega_k = 2(\beta - 1)\nabla f(x_k)p_k/(\gamma \|p_k\|^2)$$

**Proof:** Line-search will terminate as soon as $\alpha^{(\ell)} \leq \omega_k$:

1. May be that $\alpha_{\text{init}}$ satisfies the Armijo condition $\Rightarrow \alpha_k = \alpha_{\text{init}}$.

2. Otherwise in the last line-search iteration we have

$$\alpha^{(\ell-1)} > \omega_k, \qquad \alpha_k = \alpha^{(\ell)} = \tau\alpha^{(\ell-1)} > \tau\omega_k.$$

Combining these 2 cases gives the required result. $\square$

Notes

# Backtracking-Armijo line-search

1. The previous analysis permit to say that Backtracking-Armijo line-search ends in a finite number of steps.

2. The line-search produce a step length **not too long** due to the condition

$$f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \leq f(\boldsymbol{x}_k) + \alpha_k \beta \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$

3. The line-search produce a step length **not too short** due to the **finite termination** theorem.

4. Armijo line-search can be improved by adding some further requirements on the step length acceptance criteria.

Notes

# Global convergence

## Theorem (Global convergence)

*Suppose that* $f(\boldsymbol{x})$ *satisfy the standard assumptions, then, for the iterates generated by the* Generic minimization algorithm *with* backtracking Armijo line-search *either:*

1. $\nabla f(\boldsymbol{x}_k)^T = \boldsymbol{0}$ *for some* $k \geq 0$;

2. *or* $\lim_{k\to\infty} f(\boldsymbol{x}_k) = -\infty$;

3. *or* $\lim_{k\to\infty} |\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k| \min\left\{1, \|\boldsymbol{p}_k\|^{-1}\right\} = 0.$

## Remark

*If the theorem, point 1 means that we found a stationary point in a finite number of steps. Point 2 means that function* $f(\boldsymbol{x})$ *is unbounded below, so that a minimum does not exists. Point 3 alone do not imply convergence, but if* $\nabla f(\boldsymbol{x}_k)$ *and* $\boldsymbol{p}_k$ *do not become orthogonal and* $\|\boldsymbol{p}_k\| \nrightarrow 0$ *then* $\|\nabla f(\boldsymbol{x}_k)\| \to 0$.

Notes

**Proof:**(Proof. $(1/3)$) Assume points 1 and 2 are not satisfied, then we prove point 3. Consider

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) + \alpha_k \beta \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k \leq f(\boldsymbol{x}_0) + \sum_{j=0}^{k} \alpha_j \beta \nabla f(\boldsymbol{x}_j) \boldsymbol{p}_j$$

by the fact that $\boldsymbol{p}_k$ is a descent direction we have that the series:

$$\sum_{j=0}^{\infty} \alpha_j \left| \nabla f(\boldsymbol{x}_j) \boldsymbol{p}_j \right| \leq \beta^{-1} \lim_{k \to \infty} \left[ f(\boldsymbol{x}_0) - f(\boldsymbol{x}_{k+1}) \right] < \infty$$

and then

$$\lim_{j \to \infty} \alpha_j \left| \nabla f(\boldsymbol{x}_j) \boldsymbol{p}_j \right| = 0$$

Notes

**Proof:** (Proof. $(2/3)$) Recall that from finite termination Armijo theorem (slide n.28)

$$\alpha_k \geq \min\left\{\alpha_{\text{init}}, \tau\omega_k\right\}, \qquad \omega_k = 2(\beta-1)\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k/(\gamma\left\|\boldsymbol{p}_k\right\|^2)$$

and consider the two index set:

$$\mathcal{K}_1 = \left\{k \mid \alpha_k \geq \alpha_{\text{init}}\right\}, \qquad \mathcal{K}_2 = \left\{k \mid \alpha_k < \alpha_{\text{init}}\right\},$$

Obviously $\mathbb{N} = \mathcal{K}_1 \cup \mathcal{K}_2$ and from $\lim_{k\to\infty} \alpha_k\left|\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k\right| = 0$ we have

$$\lim_{k\in\mathcal{K}_1\to\infty} \alpha_k\left|\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k\right| = 0, \tag{A}$$

$$\lim_{k\in\mathcal{K}_2\to\infty} \alpha_k\left|\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k\right| = 0, \tag{B}$$

Notes

**Proof:** (Proof. $(3/3)$) For $k \in \mathcal{K}_1$ we have $\alpha_k = \alpha_{\text{init}}$ and $\alpha_k \left| \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k \right| = \alpha_{\text{init}} \left| \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k \right|$ and from $(A)$ we have

$$\lim_{k \in \mathcal{K}_1 \to \infty} \left| \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k \right| = 0 \tag{$\star$}$$

For $k \in \mathcal{K}_2$ we have $\tau \omega_k \leq \alpha_k \leq \omega_k$ so

$$\alpha_k \left| \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k \right| \geq \tau \omega_k \left| \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k \right| \geq 2\tau(1-\beta) \frac{\left| \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k \right|^2}{\gamma \left\| \boldsymbol{p}_k \right\|^2}$$

and from $(B)$ we have

$$\lim_{k \in \mathcal{K}_1 \to \infty} \frac{\left| \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k \right|}{\left\| \boldsymbol{p}_k \right\|} = 0 \tag{$\star\star$}$$

Combining $(\star)$ and $(\star\star)$ gives the required result. $\square$

Notes

# Steepest descent algorithm

## Steepest descent algorithm

Given an initial guess $x_0$, let $k = 0$;
**while** not converged **do**
    Compute a step-size $\alpha_k$ using a line-search along $-\nabla f(x_k)^T$.
    Set $x_{k+1} = x_k - \alpha_k \nabla f(x_k)^T$ and increase $k$ by 1.
**end while**

- The steepest descent algorithm is simply the **generic minimization algorithm** with search direction the opposite of the gradient in $x_k$.

- The search direction $-\nabla f(x_k)^T$ is always a **descent direction** unless the point $x_k$ is a stationary point.

Notes

# Global convergence of steepest descent

---

## Corollary (Global convergence of steepest descent)

*Suppose that $f(x)$ satisfy the standard assumptions, then, for the iterates generated by the **steepest descent algorithm** with **backtracking Armijo line-search** either:*

1. $\nabla f(x_k)^T = 0$ *for some $k \geq 0$;*
2. *or $\lim_{k \to \infty} f(x_k) = -\infty$;*
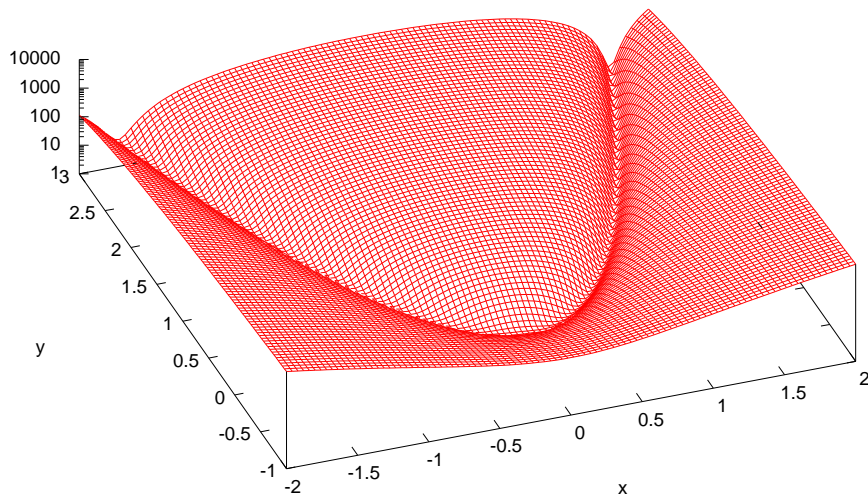3. *or $\lim_{k \to \infty} \nabla f(x_k)^T = 0$.*
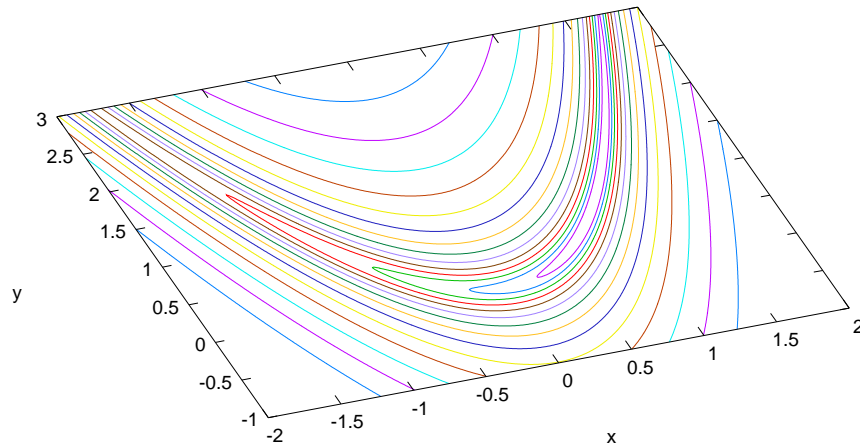
Notes

# The Rosenbrock example

■ Although the **steepest descent** scheme is globally convergent it can be very slow!

■ A classical example is the Rosenbrock function:

$$f(x,y) = 100 \, (y - x^2)^2 + (x - 1)^2$$

Notes

■ This function has a unique minimum at $(1,1)^T$ inside a **banana shaped** valley.
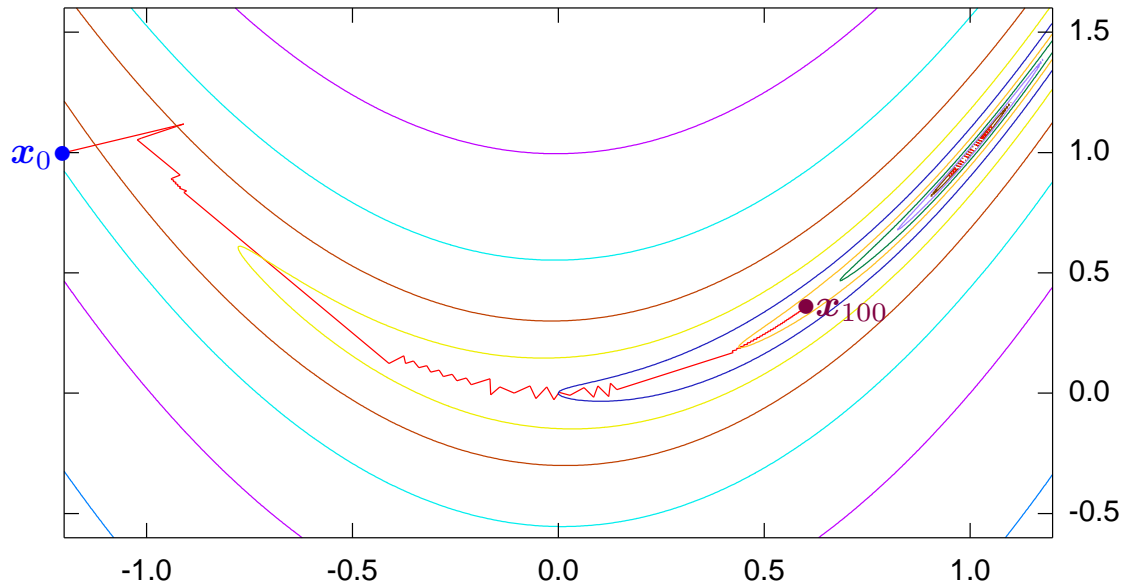
Notes

# The Rosenbrock example $(3/3)$
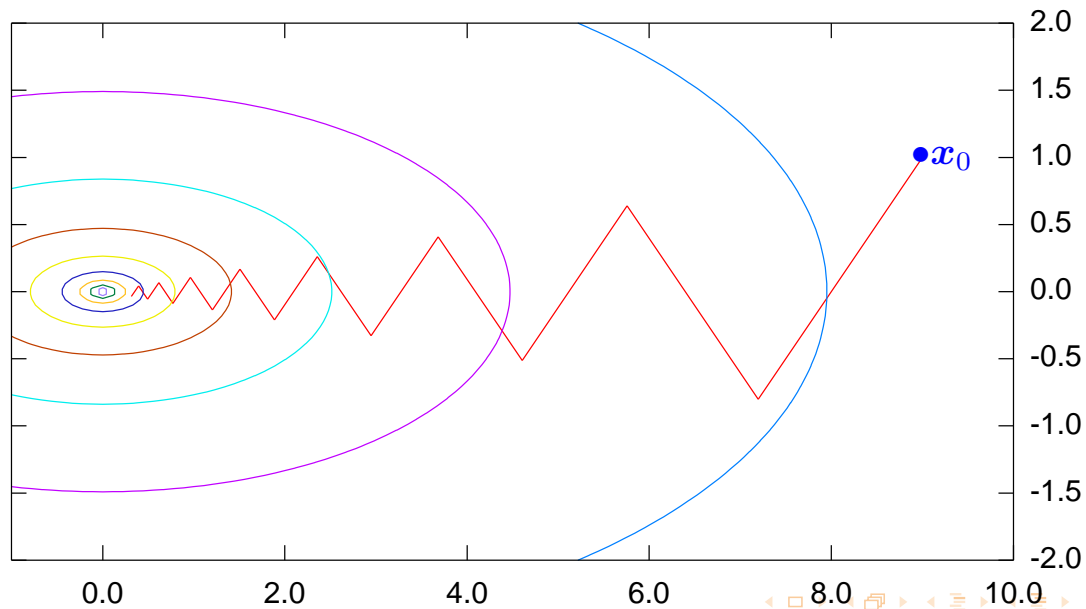
- After $100$ iteration starting from $(-1.2, 1)^T$ the approximate minimum is far from the solution.

Notes

- The steepest descent is a slow method, not only on a difficult test case like the Rosenbrock example.
- Given the function $f(x, y) = \dfrac{1}{2}x^2 + \dfrac{9}{2}y^2$ starting from $x_0 = (9, 1)^T$ we have the zig-zag pattern toward $(0, 0)^T$.

Notes

# Outline

Notes

# The Wolfe and Armijo Goldstein conditions

1. The simple condition of **descent step** is in general not enough for the convergence of a iterative minimization scheme.

2. The condition of **sufficient decrease** of backtracking Armijo line-search may be insufficient on general inexact line-search algorithm.

3. Adding another condition to the **sufficient decrease** condition such that we avoid **too short** step length we obtain **globally convergent** numerical procedure.

4. Depending on which additional condition is added we obtain the:
   1. Wolfe conditions;
   2. Armijo Goldstein conditions.

Notes

## The Wolfe conditions

Let $c_1$ and $c_2$ two constant such that $0 < c_1 < c_2 < 1$. We say that the step length $\alpha_k$ satisfy the Wolfe conditions if $\alpha_k$ satisfy:

1. sufficient decrease: $f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \leq f(\boldsymbol{x}_k) + c_1 \alpha_k \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$;

2. curvature condition: $\nabla f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \boldsymbol{p}_k \geq c_2 \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$.

$$f(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k)$$

$$f(\boldsymbol{x}_k) + \alpha c_1 \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$
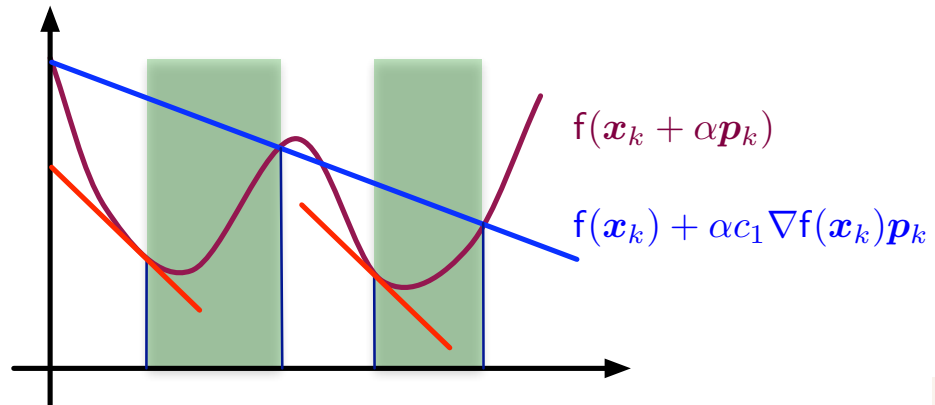
Notes

## The strong Wolfe conditions

Let $c_1$ and $c_2$ two constant such that $0 < c_1 < c_2 < 1$. We say that the step length $\alpha_k$ satisfy the strong Wolfe conditions if $\alpha_k$ satisfy:

1. **sufficient decrease:** $\mathsf{f}(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \leq \mathsf{f}(\boldsymbol{x}_k) + c_1 \, \alpha_k \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k$;

2. **curvature condition:** $|\nabla \mathsf{f}(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \boldsymbol{p}_k| \leq c_2 \, |\nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k|$.

$\mathsf{f}(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k)$

$\mathsf{f}(\boldsymbol{x}_k) + \alpha c_1 \nabla \mathsf{f}(\boldsymbol{x}_k) \boldsymbol{p}_k$

Notes

# Existence of "Wolfe" step length

- The Wolfe condition seems quite restrictive.
- The next lemma answer to the question if a step length satisfying Wolfe conditions does exists.

---

## Lemma (strong Wolfe step length)

*Let* $f : \mathbb{R}^n \mapsto \mathbb{R}$ *satisfying the regularity assumption. If the following condition are satisfied:*

1. $p_k$ *is a descent direction for the point* $x_k$, *i.e.* $\nabla f(x_k) p_k < 0$;
2. $f(x_k + \alpha p_k)$ *is bounded from below, i.e.* $\lim_{\alpha \to \infty} f(x_k + \alpha p_k) > -\infty.$

*then for any* $0 < c_1 < c_2 < 1$ *there exists an interval* $[a, b]$ *such that all* $\alpha_k \in [a, b]$ *satisfy the strong Wolfe conditions.*

---

Notes

**Proof:** Define $\ell(\alpha) = f(\boldsymbol{x}_k) + \alpha c_1 \nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k$ and $g(\alpha) = f(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k)$. From $\lim_{\alpha \to \infty} \ell(\alpha) = -\infty$ and from condition 1 it follows that there exists $\alpha_\star > 0$ such that

$$\ell(\alpha_\star) = g(\alpha_\star) \qquad \text{and} \qquad \ell(\alpha) > g(\alpha), \quad \forall \alpha \in (0, \alpha_\star)$$

so that all step length $\alpha \in (0, \alpha_\star)$ satisfy strong Wolfe condition 1. Because $\ell(0) = g(0)$ form Cauchy-Rolle theorem there exists $\alpha_{\star\star} \in (0, \alpha_\star)$ such that

$$g'(\alpha_{\star\star}) = \ell'(\alpha_{\star\star}) \qquad \Rightarrow$$

$$0 > \nabla f(\boldsymbol{x}_k + \alpha_{\star\star}\boldsymbol{p}_k)\boldsymbol{p}_k = c_1 \nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k > c_2 \nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k$$

by continuity we find an interval around $\alpha_{\star\star}$ with step lengths satisfying strong Wolfe conditions. $\square$

Notes

# The Zoutendijk condition

## Theorem (Zoutendijk)

*Let* $f : \mathbb{R}^n \mapsto \mathbb{R}$ *satisfying the regularity assumption and bounded from below, i.e.*

$$\inf_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}) > -\infty$$

*Let* $\{\boldsymbol{x}_k\}$, $k = 0, 1, \ldots, \infty$ *generated by a generic minimization algorithm where line-search satisfy Wolfe conditions, then*

$$\sum_{k=1}^{\infty} (\cos \theta_k)^2 \left\| \nabla f(\boldsymbol{x}_k)^T \right\|^2 < +\infty$$

*where*

$$\cos \theta_k = \frac{-\nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k}{\left\| \nabla f(\boldsymbol{x}_k)^T \right\| \left\| \boldsymbol{p}_k \right\|}$$

Notes

**Proof:** (Proof. $(1/3)$) Using the second condition of Wolfe (with $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k$)

$$\nabla f(\boldsymbol{x}_{k+1})\boldsymbol{p}_k \geq c_2 \nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k$$

$$\big(\nabla f(\boldsymbol{x}_{k+1}) - \nabla f(\boldsymbol{x}_k)\big)\boldsymbol{p}_k \geq (c_2 - 1)\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k$$

by using Lipschitz regularity

$$\big\|\nabla f(\boldsymbol{x}_{k+1}) - \nabla f(\boldsymbol{x}_k))\boldsymbol{p}_k\big\| \leq \gamma \|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\| \|\boldsymbol{p}_k\|$$

$$= \alpha_k \gamma \|\boldsymbol{p}_k\|^2$$

and using both inequality we obtain the estimate for $\alpha_k$:

$$\alpha_k \geq \frac{c_2 - 1}{\gamma \|\boldsymbol{p}_k\|^2} \nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k$$

Notes

**Proof:** (Proof. $(2/3)$) Using the first condition of Wolfe and lower bound estimate of $\alpha_k$

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) + \alpha_k c_1 \nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k$$

$$\leq f(\boldsymbol{x}_k) - \frac{c_1(1-c_2)}{\gamma \|\boldsymbol{p}_k\|^2}\left(\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k\right)^2$$

setting $A = c_1(1-c_2)/\gamma$ and using the definition of $\cos\theta_k$

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k) - A(\cos\theta_k)^2 \left\|\nabla f(\boldsymbol{x}_k)^T\right\|^2$$

and by induction

$$f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_1) - A\sum_{j=1}^{k}(\cos\theta_j)^2 \left\|\nabla f(\boldsymbol{x}_j)^T\right\|^2$$

Notes

**Proof:**(Proof.  (3/3)) The function $f(\boldsymbol{x})$ is bounded from below, i.e.

$$\inf_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}) > -\infty$$

so that

$$A \sum_{j=1}^{k} (\cos \theta_j)^2 \left\| \nabla f(\boldsymbol{x}_j)^T \right\|^2 \le f(\boldsymbol{x}_1) - f(\boldsymbol{x}_{k+1})$$

and

$$A \sum_{j=1}^{\infty} (\cos \theta_j)^2 \left\| \nabla f(\boldsymbol{x}_j)^T \right\|^2 \le f(\boldsymbol{x}_1) - \lim_{k \to \infty} f(\boldsymbol{x}_{k+1}) < +\infty$$

$\square$

Notes

## Corollary (Zoutendijk condition)

*Let* $f : \mathbb{R}^n \mapsto \mathbb{R}$ *satisfying the regularity assumption and bounded from below. Let* $\{\boldsymbol{x}_k\}$, $k = 0, 1, \ldots, \infty$ *generated by a generic minimization algorithm where line-search satisfy Wolfe conditions, then*

$$\cos \theta_k \left\| \nabla f(\boldsymbol{x}_k)^T \right\| \to 0 \qquad \text{where} \qquad \cos \theta_k = \frac{-\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k}{\left\| \nabla f(\boldsymbol{x}_k)^T \right\| \left\| \boldsymbol{p}_k \right\|}$$

## Remark

*If* $\cos \theta_k \geq \delta > 0$ *for all* $k$ *from the Zoutendijk condition we have:*

$$\left\| \nabla f(\boldsymbol{x}_k)^T \right\| \to 0$$

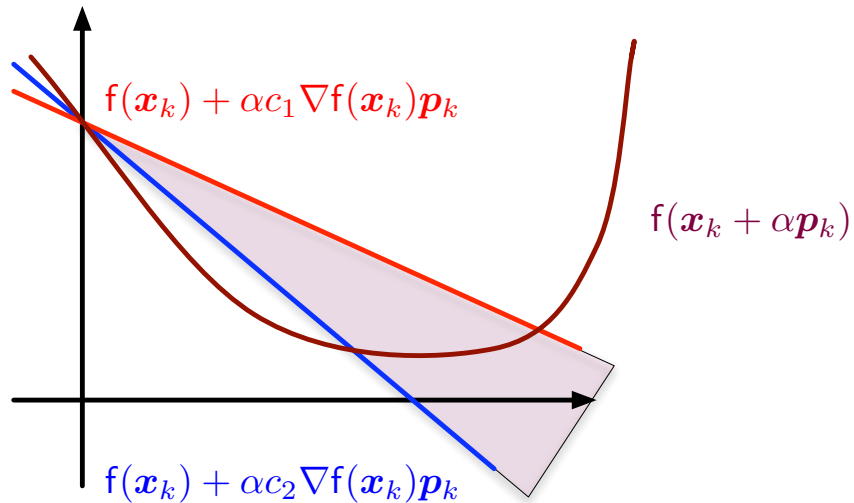*i.e. the generic minimization algorithm where line-search satisfy Wolfe conditions converge to a stationary point.*

Notes

## The Armijo-Goldstein conditions

Let $c_1$ and $c_2$ two constant such that $0 < c_1 < c_2 < 1$. We say that the step length $\alpha_k$ satisfy the Wolfe conditions if $\alpha_k$ satisfy:

1. $f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \leq f(\boldsymbol{x}_k) + c_1 \, \alpha_k \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$;

2. $f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) \geq f(\boldsymbol{x}_k) + c_2 \, \alpha_k \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$;

$$f(\boldsymbol{x}_k) + \alpha c_1 \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$

$$f(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k)$$

$$f(\boldsymbol{x}_k) + \alpha c_2 \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$

Notes

# The Armijo-Goldstein conditions

1. Armijo-Goldstein conditions has very similar theoretical properties like the Wolfe conditions.

2. Global convergence theorems can be established.

3. The weakness of Armijo-Goldstein conditions respect to Wolfe conditions is that the former can exclude local minima's from the step length as you can see in the figure below.

$$f(\boldsymbol{x}_k) + \alpha c_1 \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$

$$f(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k)$$

$$f(\boldsymbol{x}_k) + \alpha c_2 \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k$$

# Outline

Notes

# Armijo Parabolic-Cubic search

1. Backtracking-Armijo line-search can be slow if a large number of reduction must be performed to satisfy Armijo condition.

2. A better performance is obtained if instead of reducing by a fixed factor we use polynomial interpolation to estimate the location of the minimum.

3. Assuming that that $f(\boldsymbol{x}_k)$ and $\nabla f(\boldsymbol{x}_k)\boldsymbol{p}_k$ are known at the first step we know also $f(\boldsymbol{x}_k + \lambda \boldsymbol{p}_k)$ if $\lambda$ is the first trial step.

4. In this case a parabolic interpolation can be used to estimate the minimum.

5. If we store the last trial step length, in the successive iteration we can use cubic interpolation to estimate the minima's.

6. The resulting algorithm is in the following slides.

Notes

## Algorithm (Armijo Parabolic-Cubic search)

1: *armijo_linesearch*$(f, \boldsymbol{x}, \boldsymbol{p}, \tau)$
2: $f_0 \leftarrow f(\boldsymbol{x}); \ \nabla f_0 \leftarrow \nabla f(\boldsymbol{x})\boldsymbol{p}; \ \lambda \leftarrow 1;$
3: **while** $\lambda \geq \lambda_{\min}$ **do**
4:      $f_\lambda \leftarrow f(\boldsymbol{x} + \lambda\boldsymbol{p});$
5:      **if** $f_\lambda \leq f_0 + \lambda\tau\nabla f_0$ **then**
6:          **return** $\lambda$ ; *successful search*
7:      **else**
8:          **if** $\lambda = 1$ **then**
9:              $\lambda_{tmp} \leftarrow \nabla f_0 / \big[2(f_0 + \nabla f_0 - f_\lambda)\big];$
10:          **else**
11:              $\lambda_{tmp} \leftarrow$ *cubic*$(f_0, \nabla f_0, f_\lambda, \lambda, f_p, \lambda_p);$
12:          **end if**
13:          $\lambda_p \leftarrow \lambda; \ f_p \leftarrow f_\lambda; \ \lambda \leftarrow$ *range*$(\lambda_{tmp}, \lambda/10, \lambda/2);$
14:      **end if**
15: **end while**
16: **return** $\lambda_{\min}$ ; *failed search*

Notes

## Algorithm (Armijo Parabolic-Cubic search)

17: *range*$(\lambda, a, b)$

18: **if** $\lambda < a$ **then**

19:     **return** $a$;

20: **else if** $\lambda > b$ **then**

21:     **return** $b$;

22: **else**

23:     **return** $\lambda$ ;

24: **end if**

Notes

## Algorithm (Armijo Parabolic-Cubic search)

25: $cubic(f_0, \nabla f_0, f_\lambda, \lambda, f_p, \lambda_p)$

26: Evaluate:

$$\begin{pmatrix} a \\ b \end{pmatrix} = \frac{1}{\lambda^2 \lambda_p^2 (\lambda - \lambda_p)} \begin{pmatrix} \lambda_p^2 & -\lambda^2 \\ -\lambda_p^3 & \lambda^3 \end{pmatrix} \begin{pmatrix} f_\lambda - f_0 - \lambda \nabla f_0 \\ f_p - f_0 - \lambda_p \nabla f_0 \end{pmatrix}$$

27: **if** $a = 0$ **then**

28:     **return** $-\nabla f_0/(2b)$;                     *cubic is a quadratic*

29: **else**

30:     $d \leftarrow b^2 - 3\,a\,\nabla f_0$;                    *discriminant*

31:     **return** $(-b + \sqrt{d})/(3a)$;              *legitimate cubic*

32: **end if**

Notes

# Wolfe linesearch

1. Wolfe linesearch is identical to the Armijo Parabolic-Cubic search, until a point satisfying the first condition is found.
2. At this point the Armijo algorithm stop while Wolfe search try to refine the search until the second condition is satisfied.
3. If the step estimated is too short then is is enlarged until it contains a minimum.
4. If the step estimated is too long it is reduced until the second condition is satisfied.

Notes

## Algorithm (Wolfe linesearch)

1: *wolfe_linesearch*$(\mathrm{f}, \boldsymbol{x}, \boldsymbol{p}, c_1, c_2)$
2: $\mathrm{f}_0 \leftarrow \mathrm{f}(\boldsymbol{x}); \ \nabla\mathrm{f}_0 \leftarrow \nabla\mathrm{f}(\boldsymbol{x})\boldsymbol{p}; \ \lambda \leftarrow 1;$
3: **while** $\lambda \geq \lambda_{\min}$ **do**
4:     $\mathrm{f}_\lambda \leftarrow \mathrm{f}(\boldsymbol{x} + \lambda\boldsymbol{p});$
5:     **if** $\mathrm{f}_\lambda \leq \mathrm{f}_0 + \lambda c_1 \nabla\mathrm{f}_0$ **then**
6:         **go to** ZOOM; *found a $\lambda$ satisfying condition 1*
7:     **else**
8:         **if** $\lambda = 1$ **then**
9:           $\lambda_{tmp} \leftarrow \nabla\mathrm{f}_0 / \big[2(\mathrm{f}_0 + \nabla\mathrm{f}_0 - \mathrm{f}_\lambda)\big];$
10:         **else**
11:           $\lambda_{tmp} \leftarrow$ *cubic*$(\mathrm{f}_0, \nabla\mathrm{f}_0, \mathrm{f}_\lambda, \lambda, \mathrm{f}_p, \lambda_p);$
12:         **end if**
13:         $\lambda_p \leftarrow \lambda; \ \mathrm{f}_p \leftarrow \mathrm{f}_\lambda; \ \lambda \leftarrow$ *range*$(\lambda_{tmp}, \lambda/10, \lambda/2);$
14:     **end if**
15: **end while**
16: **return** $\lambda_{\min}$ ; *failed search*

Notes

## Algorithm (Wolfe linesearch)

17: *ZOOM:*

18: $\nabla f_\lambda \leftarrow \nabla f(\boldsymbol{x} + \lambda \boldsymbol{p})\boldsymbol{p}$;

19: **if** $\nabla f_\lambda \geq c_2 \nabla f_0$ **then return** $\lambda$;                     *found Wolfe point!*

20: **if** $\lambda = 1$ **then**

21:         *forward search of an interval bracketing a minimum*

22:         **while** $\lambda \leq \lambda_{\max}$ **do**

23:             $\{\lambda_p, f_p\} \leftarrow \{\lambda, f_\lambda\}$;                     *save values*

24:             $\lambda \leftarrow 2\lambda$; $f_\lambda \leftarrow f(\boldsymbol{x} + \lambda \boldsymbol{p})$;

25:             **if** *not* $f_\lambda \leq f_0 + \lambda c_1 \nabla f_0$ **then**

26:                 $\{\lambda_p, f_p\} \rightleftharpoons \{\lambda, f_\lambda\}$; **go to** *REFINE*;         *swap values*

27:             **end if**

28:             $\nabla f_\lambda \leftarrow \nabla f(\boldsymbol{x} + \lambda \boldsymbol{p})\boldsymbol{p}$;

29:             **if** $\nabla f_\lambda \geq c_2 \nabla f_0$ **then return** $\lambda$;         *found Wolfe point!*

30:         **end while**

31:         **return** $\lambda_{\max}$ ; *failed search*

32: **end if**

Notes

## Algorithm (Wolfe linesearch)

33: *REFINE:*

34: $\{\lambda_{lo}, f_{lo}, \nabla f_{lo}\} \leftarrow \{\lambda, f_\lambda, \nabla f_\lambda\}; \ \Delta \leftarrow \lambda_p - \lambda_{lo};$

35: **while** $\Delta > \epsilon$ **do**

36: $\qquad \delta\lambda \leftarrow \Delta^2 \nabla f_{lo} \big/ \big[2(f_{lo} + \nabla f_{lo}\Delta - f_p)\big];$

37: $\qquad \delta\lambda \leftarrow$ *range*$(\delta\lambda, 0.2\Delta, 0.8\Delta);$

38: $\qquad \lambda \leftarrow \lambda_{lo} + \delta\lambda; \ f_\lambda \leftarrow f(x + \lambda p);$

39: $\qquad$ **if** $f_\lambda \leq f_0 + \lambda c_1 \nabla f_0$ **then**

40: $\qquad\qquad \nabla f_\lambda \leftarrow \nabla f(x + \lambda p)p;$

41: $\qquad\qquad$ **if** $\nabla f_\lambda \geq c_2 \nabla f_0$ **then return** $\lambda;$ $\qquad$ *found Wolfe point!*

42: $\qquad\qquad \{\lambda_{lo}, f_{lo}, \nabla f_{lo}\} \leftarrow \{\lambda, f_\lambda, \nabla f_\lambda\}; \ \Delta \leftarrow \Delta - \delta\lambda;$

43: $\qquad$ **else**

44: $\qquad\qquad \{\lambda_p, f_p\} \leftarrow \{\lambda, f_\lambda\}; \ \Delta \leftarrow \delta\lambda;$

45: $\qquad$ **end if**

46: **end while**

47: **return** $\lambda;$ *failed search*

Notes

# References

📄 **Jorge Nocedal, and Stephen J. Wright**
Numerical optimization
Springer, 2006

📄 **J. E. Dennis, Jr. and Robert B. Schnabel**
Numerical Methods for Unconstrained Optimization and
Nonlinear Equations
SIAM, Classics in Applied Mathematics, **16**, 1996.

📄 **J. Stoer and R. Bulirsch**
Introduction to numerical analysis
Springer-Verlag, Texts in Applied Mathematics, **12**, 2002.

Notes