

Optimization in dimension 1

- Methods of order zero (without derivatives)
- Methods of order one and above (with derivatives)

Using derivatives...

Assumptions: f is unimodal on $[a, b]$ and is **smooth** (admits as many derivatives as we want)

Suppose that x^* is a local minimum of f on $[a, b]$

Proposition 9 (Classical result - optimality conditions)

- If $x^* \in (a, b)$ then $f'(x^*) = 0$ (x^* is a critical point)
- If $x^* = a$ then $f'(x^*) \geq 0$
- If $x^* = b$ then $f'(x^*) \leq 0$

★ The second and third conditions are called **Euler inequalities**

Towards an algorithm...

★ Direct consequence of unimodality: if $a < x^* < b$ is the minimizer of f on $[a, b]$ then

$$f'(x) < 0 \text{ for } x \in [a, x^*) \quad \text{and} \quad f'(x) > 0 \text{ for } x \in (x^*, b]$$

★ Therefore, if we choose one intermediary point $a < x_n < b$ then we know the position of x^* w.r.t. x_n by looking at $f'(x_n)$

★ Note that, compared to zero-order methods, one intermediary point is enough in order to reduce the size of the search interval

Algorithm 4 (Bisection)

Initialization: $S_0 = [a_0, b_0]$, $i = 1$

Loop:

- choose $x_i = 0.5(a_{i-1} + b_{i-1})$
- compute $f'(x_i)$
 - if $f'(x_i) < 0$ then $S_i = [x_i, b]$
 - if $f'(x_i) > 0$ then $S_i = [a, x_i]$
 - if $f'(x_i) = 0$ then $x^* = x_i$ and **stop**
- replace i with $i + 1$ and continue until the desired precision is reached

★ the third option ($f'(x_i) = 0$ can (almost) never be verified numerically) when working with **fixed machine precision** for **general functions f**

Algorithm 4 (Bisection)

Initialization: $S_0 = [a_0, b_0]$, $i = 1$

Loop:

- choose $x_i = 0.5(a_{i-1} + b_{i-1})$
- compute $f'(x_i)$
 - if $f'(x_i) \leq 0$ then $S_i = [x_i, b]$
 - if $f'(x_i) > 0$ then $S_i = [a, x_i]$
 - ~~if $f'(x_i) = 0$ then $x^* = x_i$ and stop~~
- replace i with $i + 1$ and continue until the desired precision is reached

★ the third option ($f'(x_i) = 0$ can (almost) never be verified numerically) when working with **fixed machine precision** for **general functions f**

Proposition 10

The *Bisection algorithm* converges linearly with ratio 0.5.

Proof: $|S_i| = 0.5|S_{i-1}|$ therefore

$$|x^* - x_N| \leq 0.5^N(b - a).$$

- ★ Already better than the Fibonacci/Golden search algorithms.
- ★ Is there a contradiction between the optimality of their claimed optimal rate/ratio of convergence and the result stated above?

Proposition 10

The *Bisection algorithm* converges linearly with ratio 0.5.

Proof: $|S_i| = 0.5|S_{i-1}|$ therefore

$$|x^* - x_N| \leq 0.5^N(b - a).$$

- ★ Already better than the Fibonacci/Golden search algorithms.
- ★ Is there a contradiction between the optimality of their claimed optimal rate/ratio of convergence and the result stated above?

Answer: No, since the Bisection algorithm uses **information about derivatives** $f'(x_i)$ of the function f while Fibonacci/Golden search algorithms use only **the values of f** .

Proposition 10

The *Bisection algorithm* converges linearly with ratio 0.5.

Proof: $|S_i| = 0.5|S_{i-1}|$ therefore

$$|x^* - x_N| \leq 0.5^N(b - a).$$

- ★ Already better than the Fibonacci/Golden search algorithms.
- ★ Is there a contradiction between the optimality of their claimed optimal rate/ratio of convergence and the result stated above?

Answer: No, since the Bisection algorithm uses **information about derivatives** $f'(x_i)$ of the **function** f while Fibonacci/Golden search algorithms use only **the values of** f .

- ★ Bisection method can be seen as a **search for a zero of** f' . For a general function f such that $f'(a)f'(b) \leq 0$ it will converge to a **critical point of** f

Proposition 10

The *Bisection algorithm* converges linearly with ratio 0.5.

Proof: $|S_i| = 0.5|S_{i-1}|$ therefore

$$|x^* - x_N| \leq 0.5^N(b - a).$$

- ★ Already better than the Fibonacci/Golden search algorithms.
- ★ Is there a contradiction between the optimality of their claimed optimal rate/ratio of convergence and the result stated above?

Answer: No, since the Bisection algorithm uses **information about derivatives** $f'(x_i)$ of the **function** f while Fibonacci/Golden search algorithms use only **the values of** f .

- ★ Bisection method can be seen as a **search for a zero of** f' . For a general function f such that $f'(a)f'(b) \leq 0$ it will converge to a **critical point of** f
- ★ Can we reach machine precision using the bisection method? The answer is yes: **we compare the values of** f' **with** 0!

Further improvements...

- ★ all methods presented so far possess **global linear convergence** assuming that f is **unimodal**.
- ★ Can we hope for something better?

Further improvements...

- ★ all methods presented so far possess **global linear convergence** assuming that f is **unimodal**.
- ★ Can we hope for something better?

Use **curve fitting**: approximate f **locally** by a simple function with **analytically computable minimum**.

Basic ideas:

- for each iteration: a set of **working points** for which we compute the **values** and (eventually) the **derivatives**
- construct an **approximating polynomial** p
- find **analytically the minimum of p** and **update the family of working points**

First example: Newton method

★ suppose that given x we can compute $f(x)$, $f'(x)$, $f''(x)$

Algorithm 5 (Newton method in dimension one)

Initialization: Choose the starting point x_0

Step i :

- Compute $f(x_{i-1})$, $f'(x_{i-1})$, $f''(x_{i-1})$ and approximate f around x_{i-1} by its second-order Taylor expansion

$$p(x) = f(x_{i-1}) + f'(x_{i-1})(x - x_{i-1}) + \frac{1}{2}f''(x_{i-1})(x - x_{i-1})^2.$$

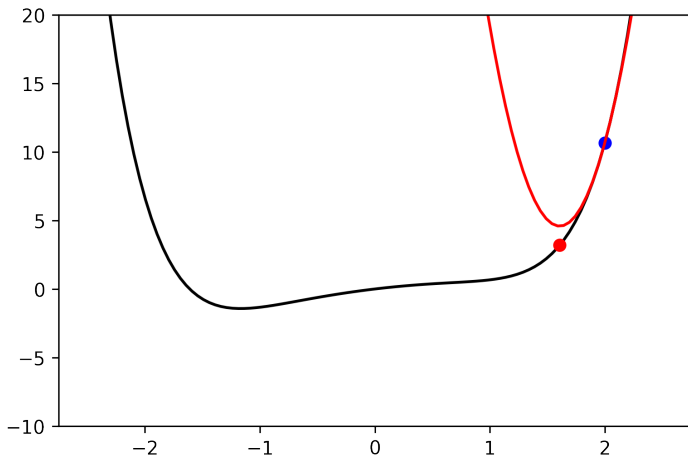
- choose x_i as the critical point of the quadratic function p :

$$x_i = x_{i-1} - \frac{f'(x_{i-1})}{f''(x_{i-1})}.$$

- replace i with $i + 1$ and loop

Example

$f(x) = x^6/6 - x^2/2 + x$ on $[-2.5, 2.5]$, $x_0 = 2$.



Proposition 11

Let $x^* \in \mathbb{R}$ be a local minimizer of a smooth function f such that $f'(x^*) = 0$ and $f''(x^*) > 0$. Then the Newton method converges to x^* *quadratically*, provided that *the starting point x_0 is close enough to x^** .

Proposition 11

Let $x^* \in \mathbb{R}$ be a local minimizer of a smooth function f such that $f'(x^*) = 0$ and $f''(x^*) > 0$. Then the Newton method converges to x^* *quadratically*, provided that *the starting point x_0 is close enough to x^** .

All the hypotheses are essential!

- What happens for $f(x) = x^4$? Which hypothesis is not verified? Does the algorithm converge for every starting point x_0 ? What is the observed convergence rate of the algorithm?
- What happens for $f(x) = \sqrt{1 + x^2}$? Does the algorithm converge for every starting point x_0 ?

Proposition 11

Let $x^* \in \mathbb{R}$ be a local minimizer of a smooth function f such that $f'(x^*) = 0$ and $f''(x^*) > 0$. Then the Newton method converges to x^* **quadratically**, provided that *the starting point x_0 is close enough to x^** .

All the hypotheses are essential!

- What happens for $f(x) = x^4$? Which hypothesis is not verified? Does the algorithm converge for every starting point x_0 ? What is the observed convergence rate of the algorithm?

Answer: $x^* = 0$, $f''(x^*) = 0$, $x_i = \frac{2}{3}x_{i-1}$. The convergence rate is **linear**.

- What happens for $f(x) = \sqrt{1+x^2}$? Does the algorithm converge for every starting point x_0 ?

Answer: $x^* = 0$, $f''(x^*) > 0$, $x_i = -x_{i-1}^3$. The convergence rate is **cubic** when $|x_0| < 1$, but the algorithm **does not converge at all** for $|x_0| \geq 1$.

- Denote $g = f'$ and observe that $g(x^*) = 0, g'(x^*) > 0$,
 $g(x^*) = g(x_i) + g'(x_i)(x^* - x_i) + \frac{1}{2}g''(\xi_i)(x^* - x_i)^2$

- Use $g(x^*) = 0$ and reformulate:

$$\frac{g(x_i)}{g'(x_i)} + (x^* - x_i) = -\frac{g''(\xi_i)}{2g'(x_i)}(x^* - x_i)^2.$$

- Use the definition of the Newton iterations to see that

$$x^* - x_{i+1} = \frac{-g''(\xi_i)}{2g'(x_i)}(x^* - x_i)^2.$$

- use the hypotheses to conclude!

Another point of view

- ★ Newton's method can be seen a linearization method for finding the zeros of $g = f'$.
- ★ Indeed, $g(x) = g(x_{i-1}) + g'(x_{i-1})(x - x_{i-1}) + o(|x - x_{i-1}|)$
- ★ Imposing that the linear part is zero amounts to

$$x = -\frac{g(x_{i-1})}{g'(x_{i-1})} + x_{i-1}$$

which is exactly the Newton method

Modified Newton: degenerate case

- ★ it is possible to show that when $f''(x^*) = 0$ then the rate of convergence is **linear**
- ★ if the multiplicity m of the root x^* of f' is known then the following modified Newton method converges quadratically (if it is well defined...)

$$x_{n+1} = x_n - m \frac{f'(x_n)}{f''(x_n)}.$$

- ★ in practice this does not really help: **you don't know the multiplicity *a priori* for a general function f !**

A second example: Regula Falsi

- ★ approximate f again by a quadratic polynomial
- ★ we consider two working points with first order information
- ★ given the two last iterates x_{i-1} and x_{i-2} we may approximate $f''(x_{i-1})$ using finite differences

$$f''(x_{i-1}) \approx \frac{f'(x_{i-1}) - f'(x_{i-2})}{x_{i-1} - x_{i-2}}$$

A second example: Regula Falsi

Algorithm 6 (False Position Method)

Initialization: Choose the starting points x_0, x_1 .

Step $i \geq 2$:

- Compute $f(x_{i-1}), f'(x_{i-1}), f'(x_{i-2})$ and approximate f around x_{i-1} with a second-order polynomial

$$p(x) = f(x_{i-1}) + f'(x_{i-1})(x - x_i) + \frac{1}{2} \frac{f'(x_{i-1}) - f'(x_{i-2})}{x_{i-1} - x_{i-2}} (x - x_{i-1})^2.$$

- choose x_i as the minimizer of the quadratic function p :

$$x_i = x_{i-1} - f'(x_{i-1}) \frac{x_{i-1} - x_{i-2}}{f'(x_{i-1}) - f'(x_{i-2})}.$$

- replace i with $i + 1$ and loop

Remarks

★ The method is symmetric with respect to x_{i-1} and x_{i-2} . It is equivalent to

$$x_i = x_{i-2} - f'(x_{i-2}) \frac{x_{i-1} - x_{i-2}}{f'(x_{i-1}) - f'(x_{i-2})}$$

★ this can be viewed again as a search for a zero of $g = f'$: approximate f' by a straight line through points $(x_{i-1}, f'(x_{i-1}))$ and $(x_{i-2}, f'(x_{i-2}))$.

★ for a non degenerate minimizer x^* of a smooth function f ($f'(x^*) = 0$, $f''(x^*) > 0$) and for x_0, x_1 close enough to x^* the method converges to x^* **superlinearly** with order of convergence

$$\lambda = (1 + \sqrt{5})/2.$$

★ the **Regula Falsi** method has a slower convergence rate than Newton's method, but it does not need the knowledge of the **second derivative**

- **Lemma:** Let (r_n) be a sequence of positive reals verifying $r_{n+1} \leq r_n r_{n-1}$ for $n \geq 1$. If $r_0, r_1 \in (0, 1)$ then
there exists a constant $C > 0$ such that $r_n \leq Cr^{\lambda^n}$,
where $r \in (0, 1)$ and $\lambda = \frac{\sqrt{5}+1}{2}$ is the golden ratio
- Show that the errors $e_n = |x^* - x_n|$ verify an inequality of the form
$$e_{n+1} \leq Me_n e_{n-1}.$$

Cubic fit

- ★ consider **two working points** x_1 and x_2 with zero and first order information
- ★ define the cubic polynomial such that

$$p(x_1) = f(x_1), p(x_2) = f(x_2), p'(x_1) = f'(x_1), p'(x_2) = f'(x_2)$$

- ★ as the next iterate, choose the local minimizer of p .
- ★ if x^* is non degenerate and the method starts **sufficiently close to** x^* then the method converges quadratically
- ★ formulas: **complicated**, if you are interested, ask for references
- ★ curve fitting is used with polynomials of small degree: **we need to be able to compute analytically position of the minima**: therefore, there is no point using **approximating polynomials of degree higher than four**!

Conclusion: curve fitting - towards descent methods

- when the algorithm works we achieve superlinear convergence
- the convergence results are local
- when applying these methods in the general case they might converge to a local maximum or a critical point
- What to do when these methods do not work?
 - alternate zero-order or bisection search methods with curve fitting (in cases where curve fitting gives iterates outside the desired search region)
 - at each iteration be sure to decrease the objective function using a line-search method

Descent direction in 1D

- if $f'(x) \neq 0$ there are only two options: go left or go right
- choose the direction $d \in \{-1, +1\}$ which decreases f .
- first order Taylor expansion:

$$f(x + \gamma d) = f(x) + \gamma d \cdot f'(x) + o(\gamma)$$

- if $d \cdot f'(x) < 0$ then if γ is small enough then

$$f(x + \gamma d) < f(x)$$

Examples when $f'(x) \neq 0$

1. $d = -f'(x)$
2. The Newton direction $d = -f'(x)/f''(x)$ is a descent direction if and only if $f''(x) > 0$.
3. The direction $d = -f'(x_{i-1}) \frac{x_{i-1} - x_{i-2}}{f'(x_{i-1}) - f'(x_{i-2})}$ from the Secant method is a descent direction if f is strictly convex.

Inexact line search

- ★ **big question**: how to choose a descent step?
- ★ the 1D reasoning will be useful in higher dimensions

Denote $q(t) = f(x + td)$ where d is a descent direction (with $d \in \{\pm 1\}$ in 1D or general in nD), sometimes called **merit function**.

- ★ Note that if d is a descent direction, then $q'(0) = d \cdot f'(x) < 0$

We perform a test for t , with three options

- a) t is good
- b) t is too big
- c) t is too small

We should be able to answer these questions by **looking at $q(t)$ and $q'(t)$** .

- ★ perform an iterative process for constructing **confidence interval $[t_l, t_r]$** for t
- ★ ideally the condition a) should be attained as **quickly as possible**!

Generic line-search algorithm

Algorithm 7 (Line-search)

Start with $t_l = 0$, $t_r = 0$ and pick an initial $t > 0$.

Iterate:

Step 1:

If a) then exit: *you found a good t*

If b) then $t_r = t$: *you found a new upper bound for t*

If c) then $t_l = t$: *you found a new lower bound for t*

Step 2:

If no valid t_r exists we choose a new $t > t_l$, like $t = 2t_l$ (extrapolation step)

Else choose a new $t \in (t_l, t_r)$, like $t = 0.5(t_l + t_r)$ (interpolation step)

- ★ a), b), c) should form a partition of \mathbb{R}_+
- ★ if t is big enough c) should be false
- ★ each interval $[t_l, t_r]$ should contain a non-trivial sub-interval verifying a)

Armijo's rule

★ $m_1 \in (0, 1)$ and $\eta > 1$ are chosen constants.

★ we fix an initial choice of $t = t_0$ (for example $t = 1$)

★ recall that $q'(0) < 0$

a) $\frac{q(t) - q(0)}{t} \leq m_1 q'(0) \iff q(t) \leq q(0) + t(m_1 q'(0))$ (t is good)

b) $m_1 q'(0) < \frac{q(t) - q(0)}{t} \iff q(t) > q(0) + t(m_1 q'(0))$ (t is too big, $t_r = t$)

c) never

★ if t is too big, then the next t is chosen as t/η (a popular choice is $\eta = 2$).

Proposition 12

Suppose that q is of class C^1 and $q'(0) < 0$. Then the line-search with Armijo's rule finishes in a finite number of steps.

Armijo's rule may lead to slow convergence: we choose once and for all a **maximal step**.

Goldstein-Price rule

★ $m_1 < m_2 \in (0, 1)$ are chosen constants

★ recall that $q'(0) < 0$

- a) $m_2 q'(0) \leq \frac{q(t) - q(0)}{t} \leq m_1 q'(0)$
 $\iff q(0) + t(m_2 q'(0)) \leq q(t) \leq q(0) + t(m_1 q'(0))$ (good t)
- b) $m_1 q'(0) < \frac{q(t) - q(0)}{t} \iff q(t) > q(0) + t(m_1 q'(0))$ (t is too big)
- c) $\frac{q(t) - q(0)}{t} < m_2 q'(0) \iff q(t) < q(0) + t(m_2 q'(0))$ (t is too small)

Proposition 13

Suppose that $q \in C^1$ is bounded from below and $q'(0) < 0$. Then the line-search with the Goldstein-Price rule finishes in a finite number of steps.

★ What about the choice of the constants m_1, m_2 ?

Wolfe rule

★ $m_1 < m_2 \in (0, 1)$ are chosen constants

★ recall that $q'(0) < 0$

a) $\frac{q(t)-q(0)}{t} \leq m_1 q'(0)$ and $q'(t) \geq m_2 q'(0)$ (good t)

b) $\frac{q(t)-q(0)}{t} > m_1 q'(0)$ (t is too big)

c) $\frac{q(t)-q(0)}{t} \leq m_1 q'(0)$ and $q'(t) < m_2 q'(0)$ (t is too small)

Proposition 14

Suppose that $q \in C^1$ is bounded from below and $q'(0) < 0$. Then the line-search with the Wolfe rule finishes in a finite number of steps.

★ The condition on $q'(t)$ is called **curvature condition**. Wolfe's rule is widely used in line-search algorithms: it gives **good convergence properties**

★ the first condition in a) assures that **the value of f decreases** while the second assures that the **slope reduces**

★ What about **the choice of the constants m_1, m_2** ?

The quadratic case

Proposition 15

Suppose that q is quadratic with minimum t^* : $q(t) = (x - t^*)^2 + a$. Then: $q'(t) = 2(x - t^*)$ and $q(t^*) = q(0) + \frac{1}{2}q'(0)t^*$.

★ we should **not refuse the optimal step** when q is quadratic!

$$\frac{q(t^*) - q(0)}{t^*} = \frac{1}{2}q'(0).$$

Armijo: $\frac{1}{2}q'(0) \leq m_1 q'(0)$

Goldstein-Price: $m_2 q'(0) \leq \frac{1}{2}q'(0) \leq m_1 q'(0)$

Wolfe: $\frac{1}{2}q'(0) \leq m_1 q'(0)$ and $q'(t^*) \geq m_2 q'(0)$

In conclusion it is recommended to:

★ choose $m_1 < 0.5$ (for Armijo, Goldstein-Price and Wolfe)

★ choose $0.5 < m_2 < 1$ (for Goldstein-Price)

Algorithm 8 (Generic gradient descent algorithm)

Initialization: Choose an initial point x_0 and the eventual parameters for the line-search algorithm

Step i :

- compute the function value $f(x_{i-1})$ and the derivative $f'(x_{i-1})$
- perform the **line-search** algorithm in order to find a **descent step t** .
- choose the next iterate

$$x_i = x_{i-1} - tf'(x_{i-1}).$$

Stopping criterion: $|f'(x_i)|$ is small, $|f(x_{i-1}) - f(x_i)|$ is small, the descent step t is too small, maximum number of iterations reached, etc.

- ★ $f'(x_{i-1})$ can be replaced with any **descent direction d** .
- ★ various simplified variants exist: fixed descent step, variable descent step
- ★ the generalization to higher dimensions is straightforward

Convergence rate?

- ★ it is a order 1 algorithm so *a priori* we cannot expect more than **linear convergence**
- ★ if $f(x) = x^2$ and we use a fixed step algorithm then the update at each iteration is

$$x_i = x_{i-1} - tf'(x_{i-1}) = (1 - 2t)x_{i-1}.$$

therefore, for $t < 0.5$ we have linear convergence to the optimum.

- ★ the function $f(x) = x^2$ is strictly convex and quadratic: the ideal case.

Therefore we cannot expect something better.

- ★ locally, around a minimizer x^* the function f is convex. Therefore, if convergence is proved for convex functions, it will follow, that locally, around the minimizer, the convergence of GD is linear

Example of global convergence result

Proposition 16 (Convergence rate for the gradient descent with fixed step)

Suppose that $f : \mathbb{R} \rightarrow \mathbb{R}$ is of class C^2 with f' Lipschitz continuous on \mathbb{R} : there exists $M > 0$ such that

$$|f'(x) - f'(y)| \leq M|x - y|, \quad \forall x, y \in \mathbb{R}.$$

Moreover, suppose that f is α -strictly convex ($f''(x) \geq \alpha > 0$) and that f is ∞ at infinity (so that a minimizer exists).

Then the Gradient Descent algorithm with fixed step t converges to the minimum linearly when t is small enough.

Proof. Define the application $\mathcal{F} : \mathbb{R} \rightarrow \mathbb{R}$

$$\mathcal{F}(x) = x - tf'(x)$$

and prove that for t small enough \mathcal{F} is a **contraction**:

$$|\mathcal{F}(x) - \mathcal{F}(y)| \leq k|x - y|, \quad k \in (0, 1).$$

★ then we know that the fixed point iteration $x_{n+1} = \mathcal{F}(x_n)$ converges to the unique fixed point, which is exactly **the optimum**.

Example of local result

Proposition 17 (Local convergence rate)

Suppose that $f : [a, b] \rightarrow \mathbb{R}$ is unimodal and has a unique minimizer x^ in $[a, b]$. Then if f is of class C^2 and $f''(x^*) > 0$ the gradient descent algorithm with fixed step t converges linearly to x^* if t is chosen small enough and x_0 is close enough to x^* .*

- ★ Taylor expansion for f' around x^* gives a recurrence relation for the error!
- ★ the condition $f''(x^*) > 0$ cannot be omitted: **degenerate minimizers will lead to sublinear rate of convergence**. Example $f(x) = x^4$.
- ★ using more involved techniques, it is possible to prove that the gradient descent **always converges to a local minimizer**, with an eventual sublinear rate of convergence
- ★ various convergence results can be formulated when using line-search procedures instead of a fixed step: **guaranteeing descent is essential for convergence**
- ★ Wolfe's rule gives good convergence results!

Improve the speed of convergence

★ we saw that Newton's method or the Secant method give **superlinear convergence** under the right hypotheses, but they **offer no guarantee of convergence**

★ modify the gradient descent algorithm by **changing the descent direction**:

$$x_{i+1} = x_i + \gamma d_i$$

where d_i is either

- $-f'(x_i)/f''(x_i)$ (if $f''(x_i) > 0$)
- $-f'(x_i) \frac{x_i - x_{i-1}}{f'(x_i) - f'(x_{i-1})}$ (if this is indeed a descent direction)

★ combine this with a line-search procedure with initial step size $t = 1$.

★ the new algorithm will **eventually attain a superlinear rate of convergence** provided we can choose the step $\gamma = 1$ for all iterations $i \geq n_0$

★ this idea is **useful in higher dimensions** where the family of descent directions is richer

Conclusions - optimization in dimension one

- there are efficient zero-order algorithms (when derivatives are not available)
- as soon as derivatives can be computed, the convergence is accelerated
- **curve-fitting** methods give increased convergence rates, but they are sensitive to the initialization
- **line-search** procedures play an important role even in higher dimensions
- **inexact line-search**: sometimes searching for an optimum **is not the main objective** but attaining a **significant decrease in the objective function** is enough
- gradient descent algorithms **(almost) always converge to a local minimizer**, but the rate of convergence is linear at best