

02610
Optimization and Data Fitting
Week 2: Line Search Methods

Yiqiu Dong

DTU Compute
Technical University of Denmark

Numerical optimization algorithms

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad f : \mathbb{R}^n \rightarrow \mathbb{R}, f \in \mathcal{C}^2(\mathbb{R}^n)$$

The iteration step in most algorithms that we will introduce in this course is essentially in the form of

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

- \mathbf{x}_k is the iterate. We expect as $k \rightarrow +\infty$ we have $\mathbf{x}_k \rightarrow \mathbf{x}^*$.
- \mathbf{p}_k is the search direction.
- α_k is the step size or step length.
- In addition, we need the stop criteria to define when the algorithm stops.

Numerical optimization algorithms

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad f : \mathbb{R}^n \rightarrow \mathbb{R}, f \in \mathcal{C}^2(\mathbb{R}^n)$$

The iteration step in most algorithms that we will introduce in this course is essentially in the form of

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

Line search Methods: The method chooses a direction \mathbf{p}_k first, then searches along this direction from the current iterate \mathbf{x}_k for a new iterate with a lower function value. For example, we can find the next step size by solving the univariate problem:

$$\min_{\alpha > 0} f(\mathbf{x}_k + \alpha \mathbf{p}_k).$$

Search directions: Steepest descent direction

Goal:

$$f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) < f(\mathbf{x}_k)$$

Applying Taylor's theorem, we have

$$f(\mathbf{x}_k + \alpha \mathbf{p}) = f(\mathbf{x}_k) + \alpha \mathbf{p}^T \nabla f(\mathbf{x}_k) + O(\alpha^2)$$

for any \mathbf{p} and α . So the search direction for *descent methods* must satisfy

$$\mathbf{p}^T \nabla f(\mathbf{x}_k) < 0.$$

A **normalized steepest descent direction** is the solution of

$$\min_{\mathbf{p}} \mathbf{p}^T \nabla f(\mathbf{x}_k), \quad \text{subject to } \|\mathbf{p}\| = 1.$$

Steepest descent methods

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

where

$$\mathbf{p}_k = \arg \min_{\mathbf{p}} \mathbf{p}^T \nabla f(\mathbf{x}_k), \quad \text{subject to } \|\mathbf{p}\| = 1.$$

If we use 2-norm on the constraint, i.e., $\|\mathbf{p}\|_2 = 1$, then we have the closed-form solution

$$\mathbf{p}_k = -\nabla f(\mathbf{x}_k) / \|\nabla f(\mathbf{x}_k)\|_2.$$

NB: The methods using $-\nabla f(\mathbf{x}_k)$ as the search direction usually are called **gradient descent methods**, but in the textbook it is called THE **steepest descent method**.

Steepest descent methods

$$\mathbf{p}_k = -\nabla f(\mathbf{x}_k) / \|\nabla f(\mathbf{x}_k)\|_2.$$

- \mathbf{p}_k is a descent direction:

$$\begin{aligned}f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) &= -\alpha_k \mathbf{p}_k^T \nabla f(\mathbf{x}_k) + O(\alpha_k^2) \\&= \alpha_k \|\nabla f(\mathbf{x}_k)\|_2 + O(\alpha_k^2)\end{aligned}$$

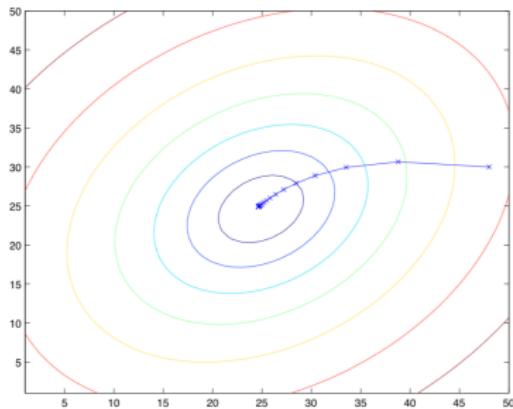
Hence, if $\nabla f(\mathbf{x}_k)_2 \neq 0$ (the first-order necessary condition is not met) and α_k is sufficiently small, we have $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$.

- \mathbf{p}_k is the max-rate descend direction: For any \mathbf{p} with $\|\mathbf{p}\|_2 = 1$, we have

$$\mathbf{p}^T \nabla f(\mathbf{x}_k) \geq -\|\mathbf{p}\|_2 \|\nabla f(\mathbf{x}_k)\|_2 = -\|\nabla f(\mathbf{x}_k)\|_2.$$

If we set $\mathbf{p} = \mathbf{p}_k$, we have $\mathbf{p}^T \nabla f(\mathbf{x}_k) = -\|\nabla f(\mathbf{x}_k)\|_2$.

Steepest descent methods



- The most popular methods (in continuous optimization)
- Simple and intuitive
- Work under very few assumptions (although they cannot directly handle non-differentiable objectives and constraints)
- Suitable for large-scale problems, e.g., easy to parallelize for problems with many terms in the objective
- Usually it's very slow

Search directions: Newton direction

We approximate $f(\mathbf{x}_k + \mathbf{p})$ by its second-order Taylor series, i.e.

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \mathbf{p}^T \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p} \stackrel{\text{def}}{=} m_k(\mathbf{p})$$

We minimize the quadratic function m_k instead of the original objective function f . Assume that the Hessian $\nabla^2 f(\mathbf{x}_k)$ is positive definite, we obtain the **Newton direction** as the minimizer of $m_k(\mathbf{p})$, i.e.,

$$\nabla^2 f(\mathbf{x}_k) \mathbf{p}_k^N = -\nabla f(\mathbf{x}_k) \implies \mathbf{p}_k^N = -(\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k).$$

- A new quadratic approximation will be constructed at \mathbf{x}_{k+1} again.
- **Special case:** The objective is quadratic. Then, the approximation is exact and the method returns a solution in one step.

Newton's method

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k^N$$

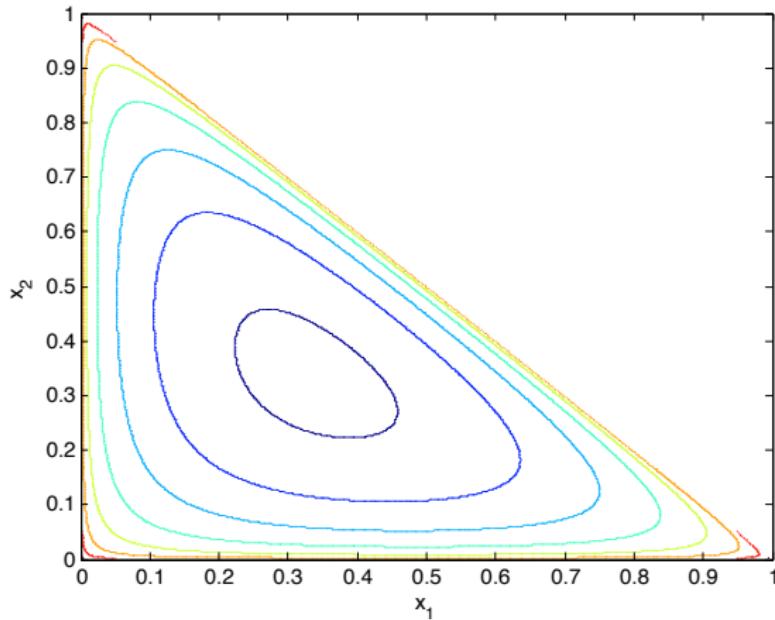
where

$$\mathbf{p}_k^N = -(\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k).$$

- Need both the gradient and the Hessian
- Based on local quadratic approximations to the objective function
- Requires a positive definite Hessian to work
- Converges very quickly near the solution
- Need a lot of computational work at each iteration:
 - ▶ Forming the Hessian
 - ▶ Inverting or factorizing the (approximate) Hessian

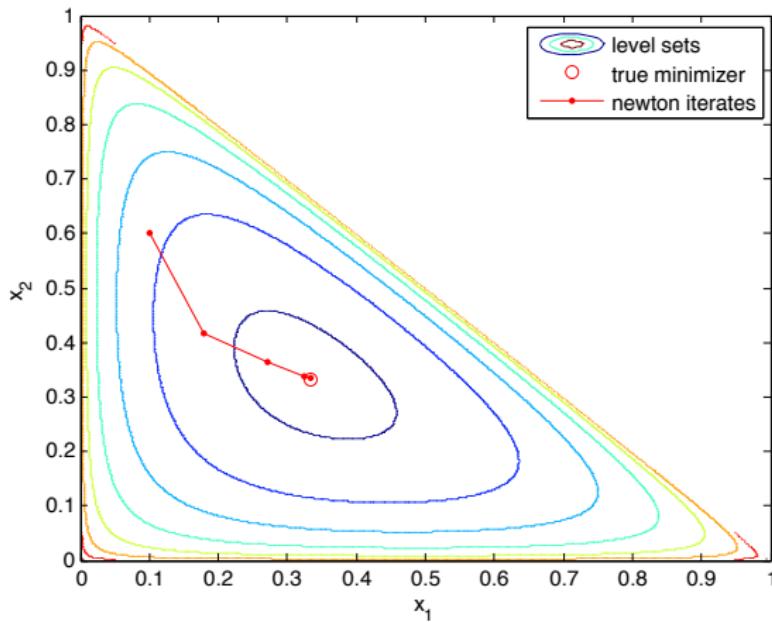
Example: Newton's method

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



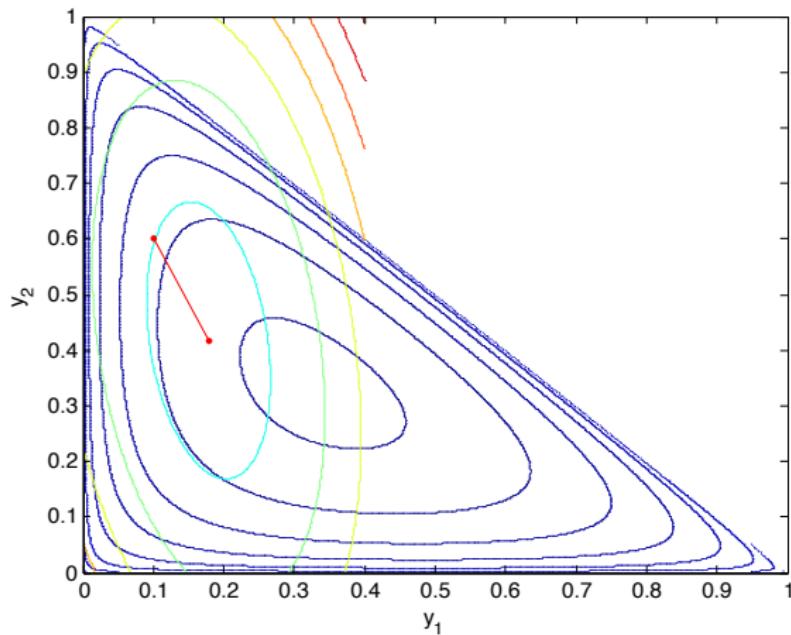
Example: Newton's method

Start the Newton's method from $x_0 = [0.1, 0.6]^T$.



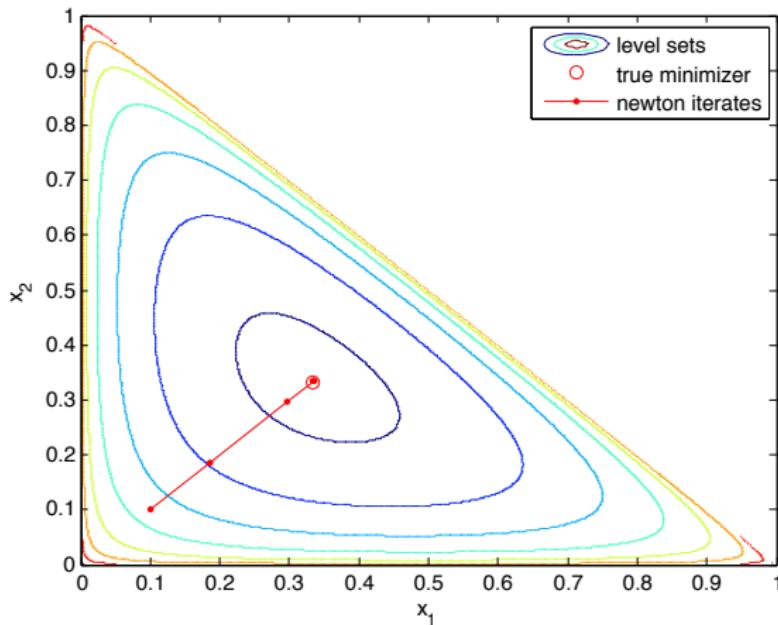
Example: Newton's method

$f(\mathbf{x}_0)$ and its quadratic approximation m_0 at $[0.1, 0.6]^T$ share the same value, gradient and Hessian. The new point minimizes m_0 .



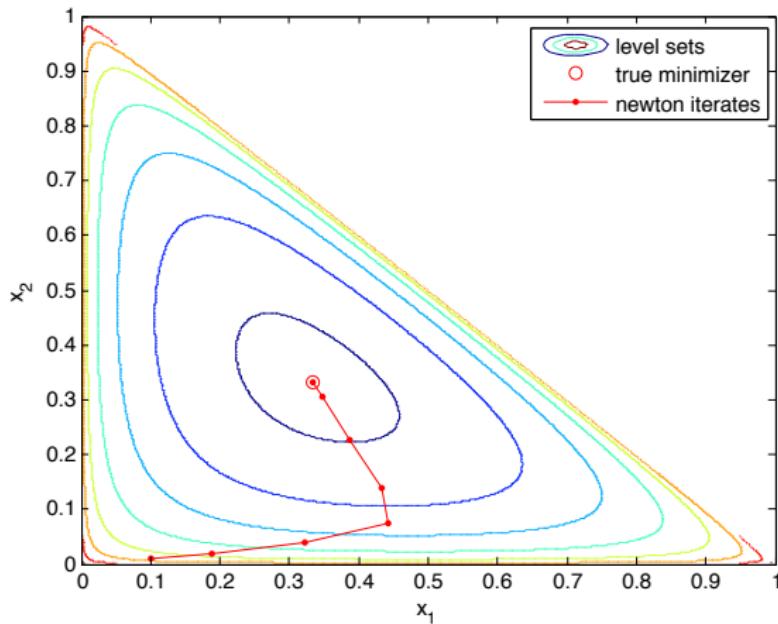
Example: Newton's method

Start the Newton's method from $x_0 = [0.1, 0.1]^T$.



Example: Newton's method

Start the Newton's method from $x_0 = [0.1, 0.01]^T$.



Newton's method: Main issues with Hessian

- **Hessian evaluation:** When the dimension n is large, obtain the Hessian can be computationally expensive.
 - ▶ **Solution:** We can use quasi-Newton methods to alleviate this difficulty. (Chapter 6)
- **Indefinite Hessian:** When the Hessian is not positive definite, the direction is not necessarily descending.
 - ▶ **Solution:** There are simple modifications. (Chapter 3.4)

Newton's method with Hessian modification

Newton direction \mathbf{p}_k^N defined by

$$\nabla^2 f(\mathbf{x}_k) \mathbf{p}_k^N = -\nabla f(\mathbf{x}_k)$$

Strategy:

- Use $\nabla^2 f(\mathbf{x}_k)$, if $\nabla^2 f(\mathbf{x}_k) \succ 0$ and its smallest eigenvalue $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) > \epsilon$; otherwise,
- use $B_k = \nabla^2 f(\mathbf{x}_k) + E_k$, so that $B_k \succ 0$ and $\lambda_{\min}(B_k) > \epsilon$.

Newton's method with Hessian modification

- **Method 1: Eigenvalue modification:** Replace negative eigenvalues by
 - ▶ a small positive number δ that is somewhat larger than machine precision u , say $\delta = \sqrt{u}$; or
 - ▶ its negative values.

It is computational expensive to find all negative eigenvalues.

- **Method 2: Adding a multiple of the identity:** Replace $\nabla^2 f(\mathbf{x}_k)$ by

$$B_k = \nabla^2 f(\mathbf{x}_k) + \tau I$$

with $\tau > 0$ such that $B_k \succ 0$.

It shifts every eigenvalue of $\nabla^2 f(\mathbf{x}_k)$ up by τ .

Newton's method with Hessian modification

- **Method 3: Modified Cholesky Factorization:** Any symmetric matrix $A \succ 0$ can be factored as

$$A = \tilde{L}\tilde{L}^T \quad \text{or} \quad A = LDL^T,$$

where \tilde{L} and L are lower triangular matrices, D is a positive diagonal matrix, and L has ones on its main diagonal.

Properties of the Cholesky factorization:

- ▶ Very useful in solving linear systems of equations. Reduce a system to two simple systems.
- ▶ The factorization is stable if A is positive definite, i.e., small errors in A will not cause large errors in L or D .
- ▶ The cost is $n^3/6 + O(n^2)$, roughly half of Gaussian elimination.

Newton's method with Hessian modification

- Method 3: Modified Cholesky Factorization:

$$A = LDL^T$$

If A is indefinite but still symmetric, then

- ▶ Cholesky factorization may not exist;
- ▶ Cholesky factorization can be not stable;
- ▶ D has zero or negative element(s) on its diagonal.

Solution:

- ▶ **Idea:** Modify A during the factorization such that all elements in D are sufficiently positive.
- ▶ **Tools:** **Pivoting**, i.e., permute the matrix at each step to pull the largest remaining diagonal element to the pivot position.
- ▶ **Effect:** Postpone the modification and keep it as small as possible.
- ▶ **Modification:** When no acceptable elements remain,

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} L_{11} & \\ L_{21} & I \end{bmatrix} \begin{bmatrix} D_1 & \\ D_2 & \end{bmatrix} \begin{bmatrix} L_{11}^T & L_{21}^T \\ & I \end{bmatrix}$$

replay D_2 by a positive definite matrix and complete the factorization.

Newton's method: Main issues with Hessian

- **Hessian evaluation:** When the dimension n is large, obtain the Hessian can be computationally expensive.
 - ▶ **Solution:** We can use quasi-Newton methods to alleviate this difficulty. (Chapter 6)
- **Indefinite Hessian:** When the Hessian is not positive definite, the direction is not necessarily descending.
 - ▶ **Solution:** There are simple modifications. (Chapter 3.4)

Search directions: Quasi-Newton direction

Consider the second-order Taylor expansion as an approximation of $f(\mathbf{x})$, i.e.

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + (\mathbf{x} - \mathbf{x}_k)^T \nabla f(\mathbf{x}_k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_k)^T \nabla^2 f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k)$$

Compute the gradient on \mathbf{x} , and obtain

$$\nabla^2 f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) \approx \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}_k).$$

Set $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ and $\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$. We choose a Hessian approximation B_{k+1} or an inverse Hessian approximation H_{k+1} satisfy the **secant equation**:

$$B_{k+1}\mathbf{s}_k = \mathbf{y}_k \quad \text{or} \quad H_{k+1}\mathbf{y}_k = \mathbf{s}_k.$$

The **quasi-Newton direction** is defined as

$$\mathbf{p}_k = -B_k^{-1}\nabla f(\mathbf{x}_k) \quad \text{or} \quad \mathbf{p}_k = -H_k\nabla f(\mathbf{x}_k).$$

Search directions

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad f : \mathbb{R}^n \rightarrow \mathbb{R}, f \in \mathcal{C}^2(\mathbb{R}^n)$$

The iteration step is essentially in the form of

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

and the search direction \mathbf{p}_k typically is from solving a linear system

$$G_k \mathbf{p} = -\nabla f(\mathbf{x}_k).$$

The difference of the search directions is mainly from how the Hessian is approximated:

$$G_k = I, \quad \text{Steepest descent}$$

$$G_k = \nabla^2 f(\mathbf{x}_k), \quad \text{Newton}$$

$$G_k = B_k, \quad \text{Quasi-Newton}$$

Step length

The iteration step is given by

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

After choosing a direction \mathbf{p}_k , we need find a step length α_k .

- Small step length:

- ▶ **Pros:** Iterations are more likely converge.
- ▶ **Cons:** Need more iterations and thus more computational work.

- Large step length:

- ▶ **Pros:** Better use the descent direction, and may reduce the total iterations.
- ▶ **Cons:** Can cause overshooting and zig-zags. If too large, even can lead to diverged iterations.

Line search

The iteration step is given by

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

The step length α_k is usually chosen by

- **Exact line search:** Find the global minimizer of the univariate problem:

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

- **Inexact line search:** Find a step length $\alpha > 0$ that only need satisfies certain conditions.

Exact line search: Convex quadratic problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}, \quad Q \succ 0$$

Then the univariate problem for exact line search is

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

According to the optimality condition for convex problems, the minimizer α^* should satisfies

$$\phi'(\alpha^*) = \alpha^* \mathbf{p}_k^T Q \mathbf{p}_k + \mathbf{p}_k^T Q \mathbf{x}_k - \mathbf{b}^T \mathbf{p}_k = 0.$$

We obtain

$$\alpha^* = -\frac{\nabla f^T(\mathbf{x}_k) \mathbf{p}_k}{\mathbf{p}_k^T Q \mathbf{p}_k}.$$

Exact line search

In general cases, we need solve a nonlinear univariate problem

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

then it is necessary to use an iterative procedure.

- According to the first-order necessary condition, we need find a positive root for the function $\phi'(\alpha)$. Then, we need to know and evaluate the derivative ϕ' .
- There exist also the methods that only need evaluate ϕ , e.g., Golden section search. (Not covered by the course.)

Methods for finding roots (Course 02601, week 2)

- **Bisection method:**

- ▶ One derivative evaluation at each iteration.
- ▶ Narrow search interval to exactly half each time.

- **Newton's method:**

- ▶ One derivative and one second derivative evaluations at each iteration.
- ▶ Must start near α^* .
- ▶ Has quadratic convergence.

- **Secant method:**

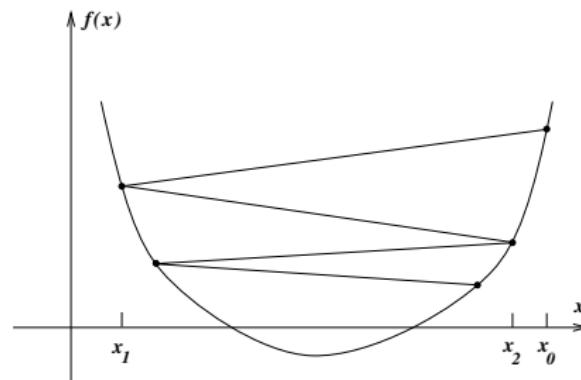
- ▶ Two points to start with; then one derivative evaluation at each iteration.
- ▶ Must start near α^* .
- ▶ Has superlinear convergence.

Inexact line search

In practice, we usually perform inexact line search to identify a step length that achieves adequate reductions in f at minimal cost.

Idea: Try out a sequence of candidate values for α , stopping to accept one when certain conditions are satisfied.

The intuitive condition $f(\mathbf{x}_k + \alpha \mathbf{p}_k) < f(\mathbf{x}_k)$ is not enough to produce convergence to \mathbf{x}^* .

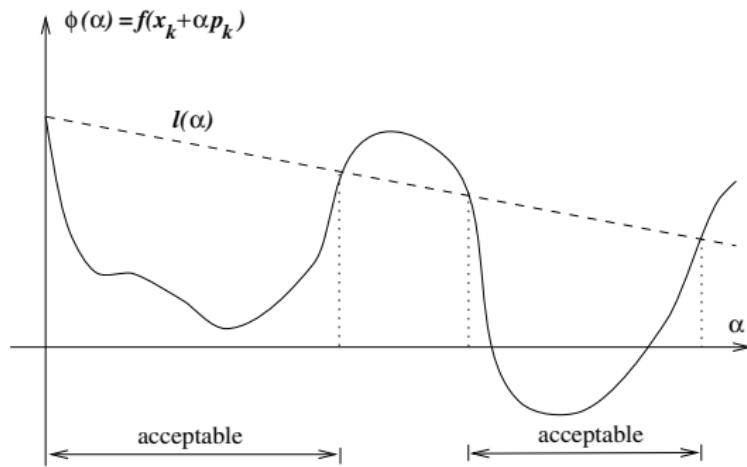


$f(\mathbf{x}^*) = -1$, but the descent sequence $\{\mathbf{x}_k\}$ gives $f(\mathbf{x}_k) = 5/k \rightarrow 0$.

Wolfe conditions

- **Armijo condition (sufficient decrease condition):**

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) \leq f(\mathbf{x}_k) + c_1 \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k \quad \text{for some } c_1 \in (0, 1).$$

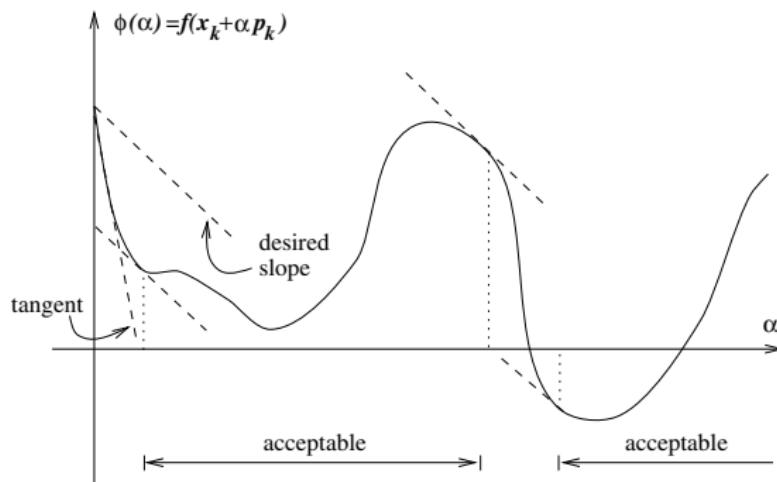


- ▶ $l(\alpha) = f(\mathbf{x}_k) + c_1 \alpha \nabla f^T(\mathbf{x}_k) \mathbf{p}_k$
- ▶ $l(0) = \phi(0)$
- ▶ The slope of $l(\alpha)$, $c_1 \nabla f^T(\mathbf{x}_k) \mathbf{p}_k$, is negative.

Wolfe conditions

- **Curvature condition:**

$$\nabla f(\mathbf{x}_k + \alpha \mathbf{p}_k)^T \mathbf{p}_k \geq c_2 \nabla f(\mathbf{x}_k)^T \mathbf{p}_k \quad \text{for some } c_2 \in (c_1, 1).$$



- ▶ Ensure the slope of ϕ at α_k is greater than c_2 times the initial slope $\phi'(0)$.
- ▶ Avoid very small α .

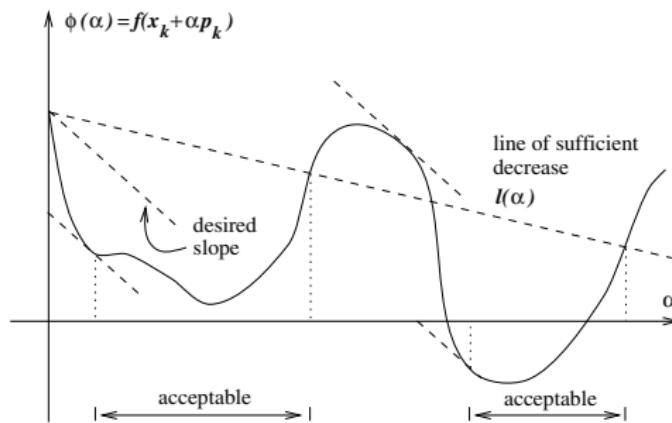
Wolfe conditions

- Wolfe conditions:

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) \leq f(\mathbf{x}_k) + c_1 \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$$

$$\nabla f(\mathbf{x}_k + \alpha \mathbf{p}_k)^T \mathbf{p}_k \geq c_2 \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$$

with $0 < c_1 < c_2 < 1$.



- The Wolfe conditions are scale-invariant.

Wolfe conditions

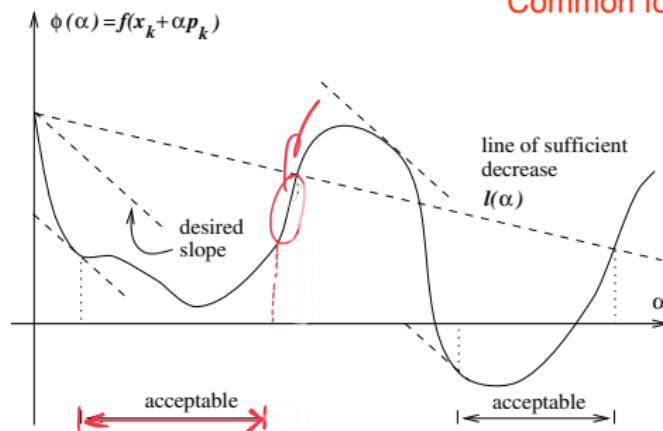
- **Strong Wolfe conditions:**

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) \leq f(\mathbf{x}_k) + c_1 \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$$

$$|\nabla f(\mathbf{x}_k + \alpha \mathbf{p}_k)^T \mathbf{p}_k| \leq c_2 |\nabla f(\mathbf{x}_k)^T \mathbf{p}_k|$$

with $0 < c_1 < c_2 < 1$.

Common for quasi-newton

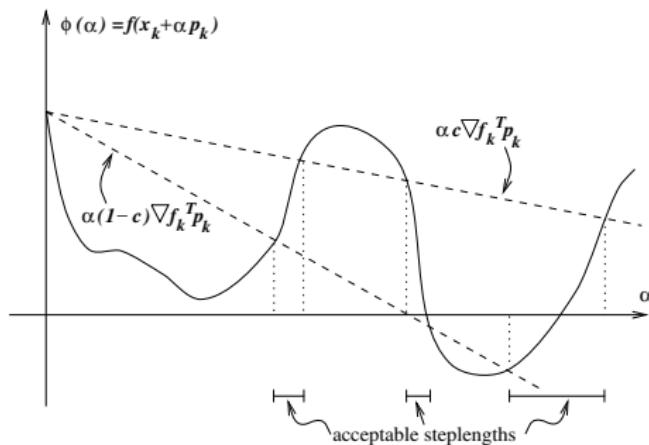


- ▶ Avoid the derivative $\phi'(\alpha_k)$ to be too positive.
- ▶ For every smooth and bounded below function f , there exist step lengths that satisfy the Wolfe condition and the strong Wolfe condition.

Goldstein conditions

$$f(\mathbf{x}_k) + (1 - c)\alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k \leq f(\mathbf{x}_k + \alpha \mathbf{p}_k) \leq f(\mathbf{x}_k) + c\alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$$

with $0 < c < \frac{1}{2}$.



- The first inequality controls the step length from below.
- The second inequality is the Armijo condition.

Backtracking line search

Idea: For a given search direction, start with a relatively large step length, then iteratively shrink the step length (i.e., “backtracking”) until a decrease of the objective function is “adequate”.

Algorithm

- ① Given a descent direction \mathbf{p}_k for f . Choose initial $\bar{\alpha} > 0$, $c \in (0, 1)$ and $\rho \in (0, 1)$.
- ② Start with $\alpha := \bar{\alpha}$
- ③ While $f(\mathbf{x}_k + \alpha \mathbf{p}_k) > f(\mathbf{x}_k) + c\alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$, set $\alpha := \rho\alpha$.
- ④ Return $\alpha_k := \alpha$

Interpolation line search

Recall $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{p}_k)$, then the Armijo condition can be written as

$$\phi(\alpha_k) \leq \phi(0) + c_1 \alpha_k \phi'(0).$$

Idea: For a given search direction, start with a relatively large step length, then iteratively shrink the step length by using the minimizer of a proper approximation of $\phi(\alpha)$ until a decrease of the objective function is "adequate".

Algorithm

- ① Given a descent direction \mathbf{p}_k for f . Choose initial $\alpha_0, c_1 \in (0, 1)$. Set $i = 0$.
- ② While $\phi(\alpha_i) > \phi(0) + c_1 \alpha_i \phi'(0)^T$,
 - ① If $i = 0$, we form a quadratic approximation to ϕ by interpolating $\phi(0)$, $\phi'(0)$ and $\phi(\alpha_0)$; otherwise, we form a cubic approximation to ϕ by interpolating $\phi(0)$, $\phi'(0)$, $\phi(\alpha_{i-1})$ and $\phi(\alpha_i)$.
 - ② Find the minimizer of the approximation to ϕ and set as α_{i+1} .
 - ③ Set $i := i + 1$.
- ③ Return $\alpha_k := \alpha_i$.

Convergence of line search methods

Theorem

Consider any iteration of the form $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$, where \mathbf{p}_k is a descent direction and α_k satisfies the Wolfe conditions. Suppose that f is bounded below in \mathbb{R}^n and that f is continuously differentiable in an open set \mathcal{N} containing the level set $\mathcal{L} \stackrel{\text{def}}{=} \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$, where \mathbf{x}_0 is the starting point of the iteration. Assume also that the gradient ∇f is Lipschitz continuous on \mathcal{N} . Then,

$$\sum_{k \geq 0} \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|_2^2 < \infty.$$

Conditions of the convergence theorem

- The angle θ_k denotes the angle between \mathbf{p}_k and the steepest descent direction $-\nabla f(\mathbf{x}_k)$, i.e,

$$\cos \theta_k = -\frac{\nabla f(\mathbf{x}_k)^T \mathbf{p}_k}{\|\nabla f(\mathbf{x}_k)\|_2 \|\mathbf{p}_k\|_2}.$$

- Lipschitz continuous gradient:** ∇f is Lipschitz continuous with constant L , if

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2$$

for all \mathbf{x} and \mathbf{y} in \mathcal{N} .

- The conditions in the theorem are not too restrictive.

Results of the convergence theorem

- The results hold also for the Goldstein conditions and strong Wolfe conditions.
- The result implies that

$$\cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|_2^2 \rightarrow 0.$$

Then, if the angle θ_k is bounded away from 90° , then we have

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0.$$

- The theorem does not guarantee the convergence to a minimizer, but only a stationary point.

Rates of convergence

Let $\{\mathbf{x}_k\}$ be a sequence in \mathbb{R}^n that converges to \mathbf{x}^* . We say that the convergence is

- **Q-linear:** if there is a constant $r \in (0, 1)$ such that

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r \quad \text{for all } k \text{ sufficiently large.}$$

- **Q-superlinear:** if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0.$$

- **Q-quadratic:** if there is a constant $M > 0$ such that

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq M \quad \text{for all } k \text{ sufficiently large.}$$

Convergence rate of steepest descent: Convex quadratic problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}, \quad Q \succ 0 \quad (1)$$

The unique minimizer \mathbf{x}^* must satisfy $Q\mathbf{x} = \mathbf{b}$.

- The search direction is $\mathbf{p}_k = -\nabla f(\mathbf{x}_k) = \mathbf{b} - Q\mathbf{x}_k$.
- The exact line search gives us the step length $\alpha_k = -\frac{\nabla f^T(\mathbf{x}_k)\mathbf{p}_k}{\mathbf{p}_k^T Q \mathbf{p}_k}$.
- The steepest descent iteration is given by $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$.

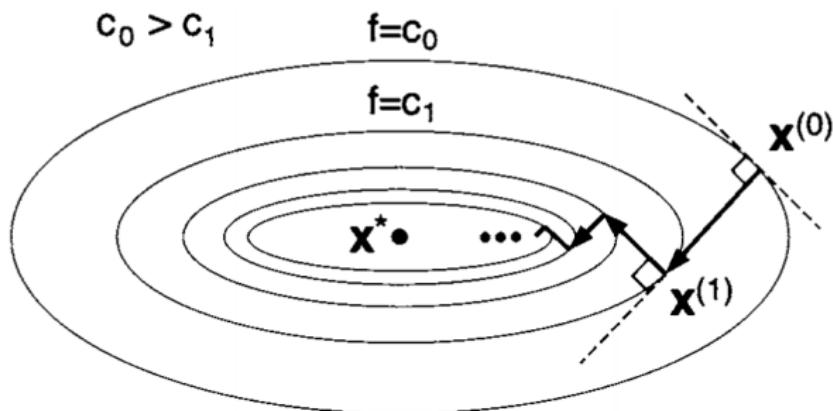
Theorem: When the steepest descent method with exact line searches is applied to (1), the error norm satisfies

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q^2 \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 \|\mathbf{x}_k - \mathbf{x}^*\|_Q^2,$$

where $\|\mathbf{x}\|_Q^2 = \mathbf{x}^T Q \mathbf{x}$, and $0 < \lambda_1 \leq \dots \leq \lambda_n$ are the eigenvalues of Q .

Examples

- **Example 1:** $f(\mathbf{x}) = x_1^2 + x_2^2$. Steepest descent arrives at $\mathbf{x}^* = \mathbf{0}$ in 1 iteration.
- **Example 2:** $f(\mathbf{x}) = \frac{1}{5}x_1^2 + x_2^2$. Steepest descent makes progress in a narrow valley



Convergence rate of steepest descent

Theorem

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable, and that the iterates generated by the steepest descent method with exact line searches converge to a point \mathbf{x}^* at which the Hessian matrix $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Let r be any scalar satisfying

$$r \in \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}, 1 \right),$$

where $\lambda_1 \leq \dots \leq \lambda_n$ are the eigenvalues of $\nabla^2 f(\mathbf{x}^*)$. Then for all k sufficiently large, we have

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq r^2(f(\mathbf{x}_k) - f(\mathbf{x}^*)).$$

- The theorem shows the **linear** convergence rate of the steepest descent methods.

Convergence rate of Newton's method

Theorem

Suppose that f is twice differentiable and that the Hessian $\nabla^2 f(\mathbf{x})$ is Lipschitz continuous in a neighborhood of a solution \mathbf{x}^* at which the sufficient conditions are satisfied. Consider the iteration $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$, where $\mathbf{p}_k = -\nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$. Then,

- ① if the starting point \mathbf{x}_0 is sufficiently close to \mathbf{x}^* , the sequence of iterates converges to \mathbf{x}^* ;
- ② the rate of convergence of $\{\mathbf{x}_k\}$ is **quadratic**; and
- ③ the sequence of gradient norms $\|\nabla f(\mathbf{x}_k)\|_2$ converges quadratically to zero.

Convergence rate of quasi-Newton method

Theorem

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable. Consider the iteration $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$ and that $\mathbf{p}_k = -B_k^{-1}\nabla f(\mathbf{x}_k)$. Let us assume also that $\{\mathbf{x}_k\}$ converges to a point \mathbf{x}^* such that $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Then $\{\mathbf{x}_k\}$ converges **superlinearly** if and only if

$$\lim_{k \rightarrow \infty} \frac{\|(B_k - \nabla^2 f(\mathbf{x}^*))\mathbf{p}_k\|_2}{\|\mathbf{p}_k\|_2} = 0.$$

- The sequence of quasi-Newton matrices B_k is not necessary to converge to $\nabla^2 f(\mathbf{x}^*)$. It suffices that the B_k become increasingly accurate approximations to $\nabla^2 f(\mathbf{x}^*)$ along the search directions \mathbf{p}_k .