# Lecture 4 Unconstrained optimization algorithms

Magnus Önnheim 2011-11-02

Consider the unconstrained optimization problem to

1

$$\underset{\mathbf{x}\in\mathbb{R}^{n}}{\operatorname{ninimize}} f(\mathbf{x}), \tag{1}$$

where  $f \in C^0$  on  $\mathbb{R}^n$  (f is continuous). Mostly, we assume that  $f \in C^1$  holds (f is continuously differentiable), often also  $f \in C^2$ 

Size of the problem (n)?

- Are  $\nabla f(\mathbf{x})$  and/or  $\nabla^2 f(\mathbf{x})$  available; to what cost?
- What it is the goal? (Global/local minimum, stationary point?)
- ▶ What are the convexity properties of *f*?
- Do we have a good estimate of the location of a stationary point x\*? (Can we use locally-only convergent methods?)

Suppose we have m data points (t<sub>i</sub>, b<sub>i</sub>) believed to be related as

$$x_1 + x_2 \exp(x_3 t_i) + x_4 \exp(x_5 t_i) = b_i, \qquad i = 1, \dots, m,$$

with unknown parameters  $x_1, \ldots, x_5$ . (Here,  $\exp(x) = e^x$ .) The best description minimizes the total "residual error," given by the norm of the residual

$$f_i(\mathbf{x}) := b_i - [x_1 + x_2 \exp(x_3 t_i) + x_4 \exp(x_5 t_i)], \quad i = 1, \dots, m$$

Resulting optimization problem:

$$\min_{\mathbf{x}\in\mathbb{R}^5} f(\mathbf{x}) := \sum_{i=1}^m |f_i(\mathbf{x})|^2 = \sum_{i=1}^m (b_i - [x_1 + x_2 \exp(x_3 t_i) + x_4 \exp(x_5 t_i)])^2$$

 Very often solved problem type within numerical analysis and mathematical statistics

#### Line search type algorithm

- Step 1. Search direction:  $\mathbf{p}_k \in \mathbb{R}^n$
- Step 2. Line search: Find  $\alpha_k > 0$  such that  $f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) < f(\mathbf{x}_k)$  holds
- Step 3. Let  $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$
- Step 4. *Termination criterion*: If fulfilled, then stop! Otherwise, let k := k + 1 and go to step 1



- The figure was plotted using several thousands of function evaluations
- Never possible in reality! (And total waste of time)
- An "orienteering map" never exists
- Most algorithms are inherently *local*, only based on info at the current point x<sub>k</sub>, that is, f(x<sub>k</sub>), ∇f(x<sub>k</sub>), and ∇<sup>2</sup>f(x<sub>k</sub>)
- Possibly also on previous points passed
- An algorithm is a "near-sighted mountain climber" when trying to reach the summit (for a max problem!)
- The mountain climber is in a deep fog and can only check her barometer for the height and feel the steepness of the slope under her feet. Go as much uphill as possible

- ► Remember descent directions are the direction where f(x<sub>k</sub> + αp) < f(x<sub>k</sub>) for α ∈ (0, δ], for some δ > 0. A few natural choices exist.
- If ∇f(x<sub>k</sub>) ≠ 0<sup>n</sup>, then p = −∇f(x<sub>k</sub>)/||∇f(x)|| is a descent direction for f at x<sub>k</sub> (Part of necessary condition proof!)
- This steepest descent direction solves the problem to

$$\min_{\mathbf{p}\in\mathbb{R}^{n}:\|\mathbf{p}\|=1} \nabla f(\mathbf{x}_{k})^{\mathrm{T}}\mathbf{p} = \min_{\mathbf{p}\in\mathbb{R}^{n}:\|\mathbf{p}\|=1} \frac{d}{d\alpha} f(\mathbf{x}_{k} + \alpha \mathbf{p})$$

► This is the "near-sighted mountain climber" direction.

Suppose Q ∈ ℝ<sup>n×n</sup> is a symmetric, positive definite matrix. Then p = −Q∇f(x<sub>k</sub>) is a descent direction for f at x<sub>k</sub>, because

$$abla f(\mathbf{x}_k)^{\mathrm{T}}\mathbf{p} = -
abla f(\mathbf{x}_k)^{\mathrm{T}}\mathbf{Q}
abla f(\mathbf{x}_k) < 0,$$

due to the positive definiteness of  ${f Q}$ 

- Special case:  $\mathbf{Q} = \mathbf{I}^n$  yields steepest descent
- Special case: Q = (∇<sup>2</sup>f(x<sub>k</sub>))<sup>-1</sup>, if the Hessian is positive definite. This is Newton's method

- Steepest descent is most often *not* a very good algorithm. (Computer exercise I)
- It fails to take into account more than information about  $\nabla f$
- A second-order Taylor approximation:

$$f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x}) \approx \varphi_{\mathbf{x}}(\mathbf{p}) := \nabla f(\mathbf{x})^{\mathrm{T}} \mathbf{p} + \frac{1}{2} \mathbf{p}^{\mathrm{T}} \nabla^{2} f(\mathbf{x}) \mathbf{p}$$

'Minimize'  $\varphi_x$  over **p** by setting the gradient of  $\varphi_x(\mathbf{p})$  to zero:

$$\nabla_{\mathbf{p}}\varphi_{\mathbf{x}}(\mathbf{p}) = \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})\mathbf{p} = \mathbf{0}^n \Longleftrightarrow \nabla^2 f(\mathbf{x})\mathbf{p} = -\nabla f(\mathbf{x})$$

• 
$$n = 1$$
:  $f'(x) + f''(x)p = 0 \implies p = -f'(x)/f''(x)$ 

▶ Provides descent if f''(x) > 0:  $f'(x)p = -[f'(x)]^2/f''(x) < 0$ 

- Methods
- ► Lack of positive definiteness. ∇<sup>2</sup>f(x) is not positive definite (PD) (that is, some eigenvalue(s) is/are negative)
- Ex: n = 1. If f"(x) < 0, then the 'minimization' of φ<sub>x</sub>(p) is actually a maximization! Bad!
- Solution: Modify  $\nabla^2 f(\mathbf{x})$  so that the result is PD.
- If λ is an eigenvalue of ∇<sup>2</sup>f(x) then for any γ ∈ ℝ, λ + γ is an eigenvalue of ∇<sup>2</sup>f(x) + γI<sup>n</sup>. Choosing γ large enough makes ∇<sup>2</sup>f(x) + γI<sup>n</sup> PD.
- ► Note: If the value of  $\gamma$  is very large  $\implies \nabla^2 f(\mathbf{x}) + \gamma \mathbf{I}^n \approx \gamma \mathbf{I}^n$  $\implies$  direction  $\approx$  steepest descent
- ► Name: *Levenberg–Marquardt*

Why do we not always choose Newton directions? II Methods

- Lack of enough differentiability. If f ∉ C<sup>2</sup>, or ∇<sup>2</sup>f(x) too impractical to compute, what do we do? Try to construct something that approximates a Hessian.
- ▶ n = 1: the *secant method*: replace f''(x) with

$$\frac{f'(x_k) - f'(x_{k-1})}{x_k - x_{k-1}}$$

▶ n > 1: *quasi-Newton*: choose approximate matrix **B**<sub>k</sub> so that

$$\mathbf{B}_k(\mathbf{x}_k - \mathbf{x}_{k-1}) = \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}),$$

+ additional choices (the above does not specify the entire matrix  $\mathbf{B}_k$ !), so that, for example,  $\mathbf{B}_{k+1}$  can be computed easily from  $\mathbf{B}_k$ , and so that  $\mathbf{B}_k$  is symmetric and positive definite

Computational burden. It may be too much to ask for to solve a linear system many times when n > 1000 or so; it is enough to do *some* work on the linear system and still get a descent property. (See The Book, p. 275 for an example)

There are many specific choices of matrices B<sub>k</sub> that lead to a variety of *quasi-Newton* methods

#### \*Additional requirements

- Purpose: prevent the descent directions to deteriorate in quality, and prevent premature convergence
- Practical criteria:

$$|\nabla f(\mathbf{x}_k)^{\mathrm{T}}\mathbf{p}_k| \ge s_1 \|\nabla f(\mathbf{x}_k)\|^2$$
, and  $\|\mathbf{p}_k\| \le s_2 \|\nabla f(\mathbf{x}_k)\|$ ,

or

$$-\frac{\nabla f(\mathbf{x}_k)^{\mathrm{T}} \mathbf{p}_k}{\|\nabla f(\mathbf{x}_k)\| \cdot \|\mathbf{p}_k\|} \ge s_1, \quad \text{and} \quad \|\mathbf{p}_k\| \ge s_2 \|\nabla f(\mathbf{x}_k)\|$$

- Interpretations: ∇f(x<sub>k</sub>)<sup>T</sup>p<sub>k</sub> is the directional derivative of f at x<sub>k</sub> in the direction of p<sub>k</sub>. Make sure it stays away from zero!
- Also, make sure that p<sub>k</sub> stays bounded and that it tends to zero if and only if ∇f(x<sub>k</sub>) does
- These conditions hold for the above examples

Approximately solve the one-dimensional problem to

$$\underset{\alpha \ge 0}{\text{minimize }} \varphi(\alpha) := f(\mathbf{x}_k + \alpha \mathbf{p}_k)$$

Its optimality conditions are that

$$\varphi'(\alpha^*) \geq \mathbf{0}, \qquad \alpha^* \cdot \varphi'(\alpha^*) = \mathbf{0}, \qquad \alpha^* \geq \mathbf{0},$$

These are the variational inequality. That is,

$$\nabla f(\mathbf{x}_k + \alpha^* \mathbf{p}_k)^{\mathrm{T}} \mathbf{p}_k \geq 0, \quad \alpha^* \cdot \nabla f(\mathbf{x}_k + \alpha^* \mathbf{p}_k)^{\mathrm{T}} \mathbf{p}_k = 0, \quad \alpha^* \geq 0,$$

holds

• If 
$$\alpha^* > 0$$
, then  $\varphi'(\alpha^*) = 0$  holds  $\Longrightarrow \nabla f(\mathbf{x}_k + \alpha^* \mathbf{p}_k)^{\mathrm{T}} \mathbf{p}_k = 0$ 

The search direction p<sub>k</sub> then is orthogonal to the gradient of f at the point x<sub>k</sub> + a<sup>\*</sup>p<sub>k</sub>



Figure: A line search in a descent direction

- No point solving the one-dimensional problem exactly! Why? The optimum to the entire problem lies elsewhere!
- Interpolation: Use f(x<sub>k</sub>), ∇f(x<sub>k</sub>)<sup>T</sup>p<sub>k</sub>, f(x<sub>k</sub> + p<sub>k</sub>) to model a quadratic function approximating f along p<sub>k</sub>. Minimize it by using the analytic formula for quadratics
- Newton's method: Repeat the improvements gained from a quadratic approximation: α := α − φ'(α)/φ"(α)
- Golden section: Derivative-free method that shrinks an interval wherein a solution to φ'(α) = 0 lies (similar to interval-halving for solving f(x) = 0)



Figure: Choosing  $\alpha$  too large may cause oscillations.

### Step length too short

## Methods



Figure: Choosing  $\alpha$  too small may stall progress.

## Armijo rule, I

- ▶ Idea: quickly generate a step  $\alpha$  which provides "sufficient" decrease in f. Note:  $f(\mathbf{x}_k + \alpha \mathbf{p}_k) \approx f(\mathbf{x}_k) + \alpha \cdot \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$ , valid for small values of  $\alpha > 0$
- ▶ Requirement: we get a decrease in f which is at least a *fraction* of that predicted in the right-hand side above. Let  $\mu \in (0, 1)$  be this fraction. Acceptable step lengths are  $\alpha > 0$  satisfying

$$\varphi(\alpha) - \varphi(0) \le \mu \alpha \varphi'(0),$$
 (2a)

that is,

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) - f(\mathbf{x}_k) \le \mu \alpha \nabla f(\mathbf{x}_k) \mathbf{p}_k$$
(2b)

In practice: Initial guess α. If (2) not satisfied, try α/2; repeat.



Figure: The interval (R) accepted by the Armijo step length rule

## \*Typical convergence result

- Suppose f ∈ C<sup>1</sup>, and for the starting point x<sub>0</sub> the level set lev<sub>f</sub> (f(x<sub>0</sub>)) = { x ∈ ℝ<sup>n</sup> | f(x) ≤ f(x<sub>0</sub>) } is bounded. Consider the iterative algorithm, with the following choices for each k:
  - **p**<sub>k</sub> satisfies the second sufficient descent condition (see \*-slide above);
  - $\|\mathbf{p}_k\| \leq M$ , where M is some positive constant; and
  - the Armijo step length rule is used

Then, the sequence  $\{\mathbf{x}_k\}$  is bounded, the sequence  $\{f(\mathbf{x}_k)\}$  is descending, lower bounded and therefore converges, and every limit point of  $\{\mathbf{x}_k\}$  is stationary

► For convex *f* much stronger convergence properties:

*Optimum exists*  $\iff$  {**x**<sub>*k*</sub>*} converges to an optimal solution* 

Theorem is a consequence of a more general theorem. Master-classes!

- Lesson 1: Cannot terminate based on the exact optimality conditions, because ∇f(x) = 0<sup>n</sup> rarely happens!
- Recommended/Common:

1. 
$$\|\nabla f(\mathbf{x}_k)\| \leq \varepsilon_1(1+|f(\mathbf{x}_k)|), \ \varepsilon_1 > 0$$
 small;

2.  $f(\mathbf{x}_{k-1}) - f(\mathbf{x}_k) \le \varepsilon_2(1 + |f(\mathbf{x}_k)|)$ ,  $\varepsilon_2 > 0$  small; and

3. 
$$\|\mathbf{x}_{k-1} - \mathbf{x}_k\| \le \varepsilon_3 (1 + \|\mathbf{x}_k\|), \ \varepsilon_3 > 0$$
 small

- Why? Need to cover cases of very steep and very flat functions
- ▶ May need to use ∞-*norm*:  $\|\mathbf{x}\|_{\infty} := \max_{1 \le j \le n} |x_j|$ , for large *n*

Problem with the scaling of the problem: If

$$\mathbf{x}_{k-1} = (56.00, 0.042, 0.000001)^{\mathrm{T}},$$
  
 $\mathbf{x}_{k} = (56.01, 0.0421, 0.000009)^{\mathrm{T}};$ 

$$\|\mathbf{x}_{k-1} - \mathbf{x}_{k}\|_{\infty} = \|(-0.01, -0.001, -0.000008)^{\mathrm{T}}\|_{\infty}$$
  
= 0.01

- Small absolute error but large relative error (800% in the third coordinate)!
- Better to apply the algorithm from a scaled problem where elements of x have similar magnitude.
- Many out-of-the-box algorithms do not check this. User responsibility!

Suppose f is only in  $C^0$ , not  $C^1$ . Example:

$$f(\mathbf{x}) := \max_{i \in \{1,...,m\}} \{ \mathbf{c}_i^{\mathrm{T}} \mathbf{x} + b_i \}, \qquad \mathbf{x} \in \mathbb{R}^n$$

- This is a piece-wise linear and convex function (see next page)
- It is differentiable almost everywhere, but not at the optimal solution!
- Ignoring non-differentiability may lead to the convergence to a non-optimal point. In other words, methods for minimizing non-differentiable function cannot only rely on gradient values
- Convex functions always has *subgradients*, corresponding to all the possible slopes of the function (Lecture 9).



Figure: A piece-wise linear convex function

Common in engineering and natural science applications that f is not explicitly given but through a simulation:

$$\mathbf{x} \in \mathbb{R}^n \longrightarrow \mathsf{Simulation} \qquad \mathbf{y} \in \mathbb{R}^m$$

- The wish is to minimize a function of both x and y: f(x, y); find the vector x that gives the best response y for f
- The form of the response y = y(x) from the input x is normally unknown
- Cannot differentiate  $\mathbf{x} \mapsto f(\mathbf{x}, \mathbf{y}(\mathbf{x}))$
- Two distinct possibilities!

- ► (1) Numerical differentiation of f by using a difference formula:
- Let  $\mathbf{e}_i = (0, 0, \dots, 0, 1, 0, \dots, 0)^{\mathrm{T}}$  be the unit vector in  $\mathbb{R}^n$ . Then,

$$f(\mathbf{x} + \alpha \mathbf{e}_i) = f(\mathbf{x}) + \alpha \mathbf{e}_i^{\mathrm{T}} \nabla f(\mathbf{x}) + (\alpha^2/2) \mathbf{e}_i^{\mathrm{T}} \nabla^2 f(\mathbf{x}) \mathbf{e}_i + \dots$$
  
=  $f(\mathbf{x}) + \alpha \partial f(\mathbf{x}) / \partial x_i + (\alpha^2/2) \partial^2 f(\mathbf{x}) / \partial x_i^2 + \dots$ 

$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + \alpha \mathbf{e}_i) - f(\mathbf{x})}{\alpha} \quad \text{(forward difference)}$$
$$\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + \alpha \mathbf{e}_i) - f(\mathbf{x} - \alpha \mathbf{e}_i)}{2\alpha} \quad \text{(central difference)}$$

- Value of α typically set to a function of the machine precision; too large → bad approximation of the partial derivative; too small → numerical cancellation
- May work well *if* the simulation is *accurate*, otherwise bad derivative information. Requires *cheap* simulations!
- ▶ (2) Derivative-free methods are available. (Not counting subgradient methods, because they demand f to be convex!) Either builds explicit models f of the objective function from evaluations of f at test points, or evaluates f at grid points that are moved around, shrunk or expanded. Names and terms: Nelder-Mead, Pattern search, surrogate models, radial basis functions, ...

- Alternative: create explicit algebraic (e.g., polynomial) model *f* based on visited points x<sub>k</sub>; solve this problem with gradient methods; evaluate its optimum in the real problem (i.e., perform a simulation); update *f* with the new information ⇒ minimizes the number of simulations!
- Recent application: diesel engine optimization for Volvo Powertrain and Volvo Car Corporation
- Optimize fuel consumption, keep soot/nitrogen emissions at an acceptable level
- Simulations hard (42 hours each) and response contains noise
- New method developed based on approximate (surrogate) models

Remember the line search algorithm form:

- Step 0. Starting point:  $\mathbf{x}_0 \in \mathbb{R}^n$ . Set k := 0Step 1. Search direction:  $\mathbf{p}_k \in \mathbb{R}^n$ Step 2. Line search: Find  $\alpha_k > 0$  such that  $f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) < f(\mathbf{x}_k)$  holds Step 3. Let  $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$
- Step 4. *Termination criterion*: If fulfilled, then stop! Otherwise, let k := k + 1 and go to step 1

Trust region does steps 1 and 2 simultaneously.

- Trust region methods often use quadratic models (as does Newton)
- Avoids line searches by bounding the length of the search direction, at the same time influencing its direction

• Let 
$$\psi_k(\mathbf{p}) := f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^{\mathrm{T}} \mathbf{p} + \frac{1}{2} \mathbf{p}^{\mathrm{T}} \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$

- The accuracy of this approximation decreases when x<sub>k</sub> + p is far away from x<sub>k</sub>
- We trust the model ψ<sub>k</sub> to be a good approximation of f(x<sub>k</sub> + p) only in a neighbourhood of x<sub>k</sub> : ||p|| ≤ Δ<sub>k</sub>
- ▶ Very useful when  $\nabla^2 f(\mathbf{x}_k)$  is not positive semi-definite

- ► Relatively easy (but not trivial) to minimize  $\psi_k(\mathbf{p})$  subject to  $\|\mathbf{p}\| \leq \Delta_k$
- Idea: when ∇<sup>2</sup>f(x<sub>k</sub>) is badly conditioned, Δ<sub>k</sub> should be small (more of a steepest descent method); if well conditioned, Δ<sub>k</sub> should be large to allow for unit steps (Newton! fast convergence)
- If  $\Delta_k$  is small enough,  $f(\mathbf{x}_k + \mathbf{p}_k) < f(\mathbf{x}_k)$  holds
- ► Even if  $\nabla f(\mathbf{x}_k) = \mathbf{0}^n$  holds,  $f(\mathbf{x}_k + \mathbf{p}_k) < f(\mathbf{x}_k)$  still holds, if  $\nabla^2 f(\mathbf{x}_k)$  is not positive definite
- Progress from stationary points if at saddle points or local maxima
- This method type is what many of MATLABs optimization routines use by default.

► Update of trust region size based on a measure of similarity between the model ψ<sub>k</sub> and f: Let

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{p}_k)}{f(\mathbf{x}_k) - \psi_k(\mathbf{p}_k)} = \frac{\text{actual reduction}}{\text{predicted reduction}}$$

If 
$$\rho_k \leq \mu$$
 let  $\mathbf{x}_{k+1} = \mathbf{x}_k$  (unsuccessful step), else  $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$  (successful step)  
Value of  $\Delta_{k+1}$  depends on  $\rho_k$ :

 $\underline{}_{k+1} = \underline{}_{k+1} = \underline{}_$ 

$$\begin{aligned}
\rho_k &\leq \mu \Longrightarrow \Delta_{k+1} = \frac{1}{2}\Delta_k, \\
\mu &< \rho_k &< \eta \Longrightarrow \Delta_{k+1} = \Delta_k, \\
\rho_k &\geq \eta \Longrightarrow \Delta_{k+1} = 2\Delta_k
\end{aligned}$$

#### \*Trust region methods, V



Figure 6 illustrates the trust region subproblem



Figure: Trust region and Newton step. The ellipses are level curves of the quadratic model; the circle defines the trust region