# Lecture 3: Unconstrained optimization algorithms

Michael Patriksson

27 January 2003

# Method of choice

Consider the unconstrained optimization problem to

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \text{imize } f(\boldsymbol{x}), \tag{1}$$

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), sometimes even $C^2$.

- Size of the problem ($n$)?
- Are $\nabla f(\boldsymbol{x})$ and/or $\nabla^2 f(\boldsymbol{x})$ 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 $\boldsymbol{x}_*$? (Can we use locally-only convergent methods?)

# Example: curve fitting by least-squares

- Suppose we have $m$ data points $(t_i, b_i)$ believed to be related as

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

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

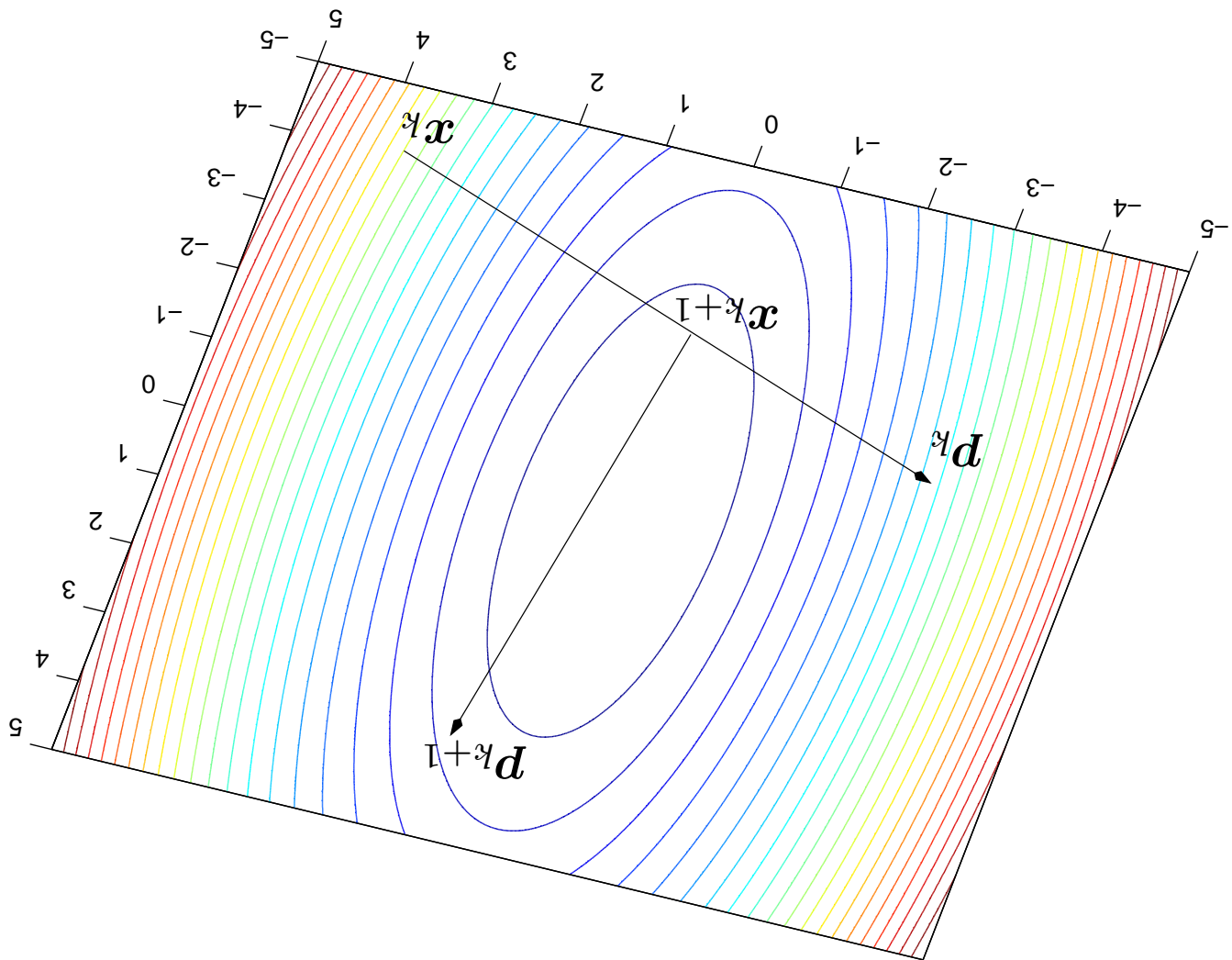$$f_i(\boldsymbol{x}) := b_i - [x_1 + x_2 \exp(x_3 t_i) + x_4 \exp(x_5 t_i)], \qquad i = 1, \ldots, m.$$

- Resulting optimization problem:

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

- Very often solved problem type within numerical analysis and mathematical statistics.

# Typical algorithm

**Step 0.** *Starting point:* $x_0 \in \mathbb{R}^n$. Set $k := 0$.

**Step 1.** *Search direction:* $p_k \in \mathbb{R}^n$.

**Step 2.** *Step length:* $\alpha_k > 0$ such that $f(x_k + \alpha_k p_k) < f(x_k)$ holds.

**Step 3.** Let $x_{k+1} := x_k + \alpha_k p_k$.

**Step 4.** *Termination criterion:* If fulfilled, then stop! Otherwise, let $k := k + 1$ and go to step 1.

# Notes

- The figure was plotted using several thousands of function evaluations.
- Never possible in reality! (And total waste of time.)
- An "orienteering map" never exists.
- Algorithms are inherently *local*, only based on info at the current point $x_k$, that is, $f(x_k)$, $\nabla f(x_k)$, and $\nabla^2 f(x_k)$.
- Possibly also on previous points passed.
- An algorithm is a "near-sighted mountain climber" when trying to reach the summit.
- 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.

# Step 1: Search directions.

- If $\nabla f(x_0) \neq \mathbf{0}^n$, then $p = -\nabla f(x_0)$ is a descent direction for $f$ at $x_0$. (Part of necessary condition proof!)

- This *steepest descent direction* solves the problem to

$$\underset{p \in \mathbb{R}^n, \|p\|=1}{\text{minimize}} \ \nabla f(x)^{\mathrm{T}} p.$$

- Suppose $Q \in \mathbb{R}^{n \times n}$ is a symmetric, positive definite matrix. Then $p = -Q\nabla f(x_0)$ is a descent direction for $f$ at $x_0$, because

$$\nabla f(x_0)^{\mathrm{T}} p = -\nabla f(x_0)^{\mathrm{T}} Q \nabla f(x_0) < 0,$$

due to the positive definiteness of $Q$.

- Special case: $Q = I^n$ yield steepest descent.

- Special case: $Q^{-1} = \nabla^2 f(x_0)$, if the Hessian is positive definite. This is *Newton's method*.

# Additional requirements

$$|\nabla f(x_k)^T p_k| \geq s_1 \|\nabla f(x_k)\|^2, \quad \text{and} \quad \|p_k\| \leq s_2 \|\nabla f(x_k)\|,$$

or

$$-\frac{\nabla f(x_k)^T p_k}{\|\nabla f(x_k)\| \cdot \|p_k\|} \geq s_1, \quad \text{and} \quad \|p_k\| \geq s_2 \|\nabla f(x_k)\|.$$

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

- $\nabla f(x_k)^T p_k$ is the directional derivative of $f$ at $x_k$ in the direction of $p_k$. Make sure it stays away from zero!

- Also, make sure that $p_k$ stays bounded and that it tends to zero if and only if $\nabla f(x_k)$ does.

- These conditions hold for the above examples.

# Newton's method

- Steepest descent is most often not a very good algorithm. Why?

- It fails to take into account more than information about $\nabla f$.

- Let

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

Minimize by setting gradient of $\varphi_{\boldsymbol{x}}(\boldsymbol{p})$ to zero:

$$\nabla_{\boldsymbol{p}}\varphi_{\boldsymbol{x}}(\boldsymbol{p}) = \nabla f(\boldsymbol{x}) + \nabla^2 f(\boldsymbol{x})\boldsymbol{p} = \boldsymbol{0}^n.$$

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

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

- Corresponding story in $\mathbb{R}^n$: $\boldsymbol{p} = -[\nabla^2 f(\boldsymbol{x})]^{-1}\nabla f(\boldsymbol{x})$, yields descent at non-stationary points if $\nabla^2 f(\boldsymbol{x})$ is positive definite!

# Why do we not always choose Newton directions?

- **Lack of positive definiteness.** $\nabla^2 f(\boldsymbol{x})$ is not positive definite (PD). Solution: add diagonal matrix so that the result is PD: $\nabla^2 f(\boldsymbol{x}) + \gamma \boldsymbol{I}^n$ for $\gamma > 0$ large enough.

- Note: If value of $\gamma$ is very large $\implies \sim$ steepest descent.

- Name: *Levenberg–Marquardt*.

- **Lack of enough differentiability.** If $f \notin C^2$, what do we do?

- $n = 1$: the secant method:

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

- $n > 1$: *quasi-Newton*: choose approximate matrix $\boldsymbol{B}_k$ so that

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

and more choices (the above does not specify the entire matrix!).

- **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 notes for an example.)

# Step 2: Line search

- Approximately solve the one-dimensional problem to

$$\underset{\alpha \geq 0}{\text{minimize }} \varphi(\alpha) := f(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k).$$

Its optimality conditions are that

$$\varphi'(\alpha_*) \geq 0, \qquad \alpha_* \cdot \varphi'(\alpha_*) = 0, \qquad \alpha_* \geq 0,$$

that is,

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

holds.

- If $\alpha_* > 0$, then $\varphi'(\alpha_*) = 0$ holds, hence $\nabla f(\boldsymbol{x}_k + \alpha_* \boldsymbol{p}_k)^{\mathrm{T}} \boldsymbol{p}_k = 0.$

- The search direction $\boldsymbol{p}_k$ is orthogonal to the gradient of $f$ at the point $\boldsymbol{x}_k + \alpha_* \boldsymbol{p}_k.$

$\varphi(\alpha)$

$\alpha^*$

$\alpha$

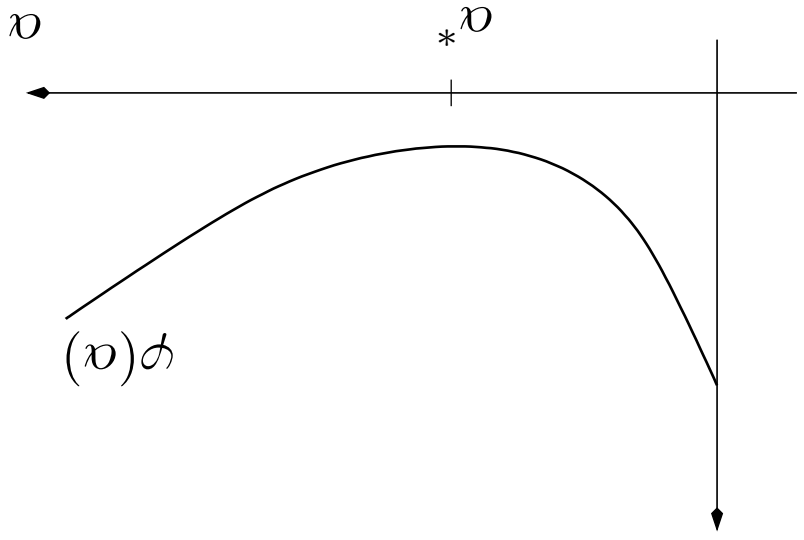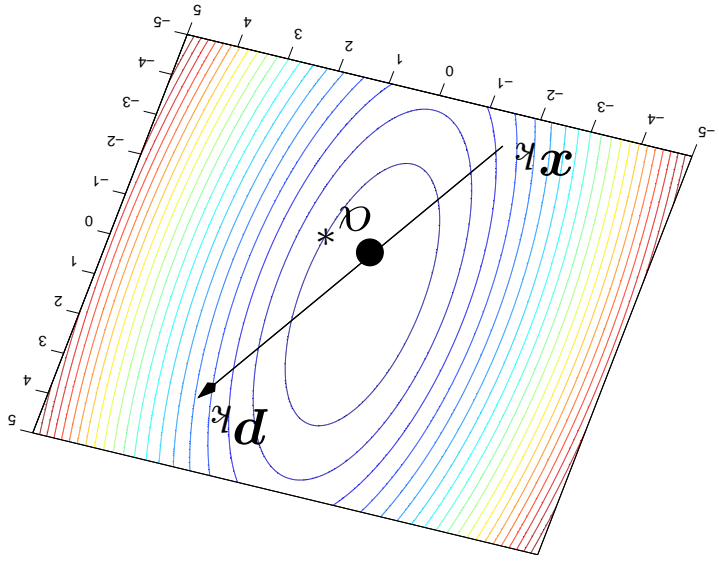$\boldsymbol{x}_k$

$\boldsymbol{d}_k$

$\alpha^*$

Figure 1: A line search in a descent direction.

# Approximate line search

- No point solving the one-dimensional problem exactly! Why? The optimum to the entire problem lies elsewhere!

- `Interpolation`: Use $f(\boldsymbol{x}_k)$, $\nabla f(\boldsymbol{x}_k)$, $\nabla f(\boldsymbol{x}_k)^{\mathrm{T}} \boldsymbol{p}_k$ to model a quadratic function approximating $f$ along $\boldsymbol{p}_k$. Minimize it by using the analytic formula for quadratics.

- `Newton's method`: Repeat the improvements gained from a quadratic approximation: $\alpha := \alpha - \varphi'(\alpha)/\varphi''(\alpha)$.

- `Golden section`: Derivative-free method that shrinks an interval where $\varphi'(\alpha) = 0$ lies.
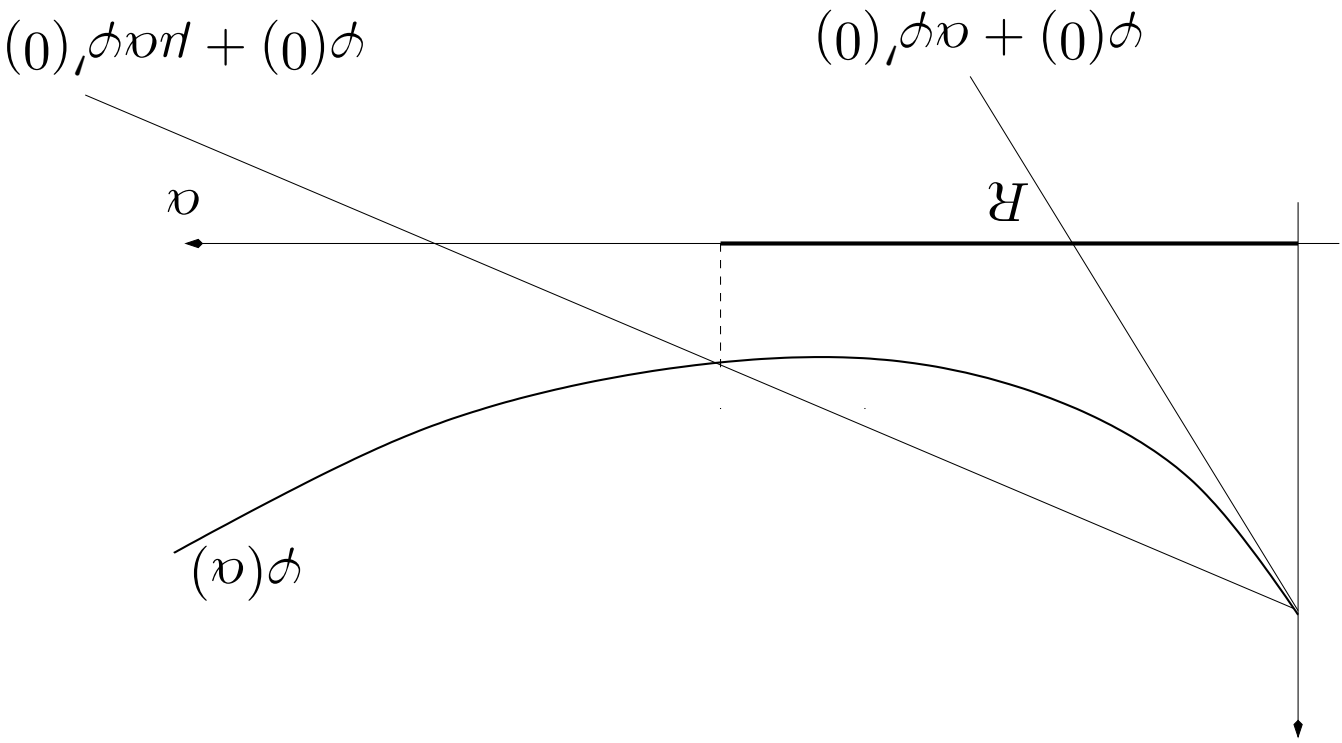
# Armijo rule

- **Idea:** quickly generate a step $\alpha$ which provides "sufficient" decrease in $f$. Note: $f(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k) \approx f(\boldsymbol{x}_k) + \alpha \cdot \nabla f(\boldsymbol{x}_k)^{\mathrm{T}} \boldsymbol{p}_k$, 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

$$\phi(\alpha) - \phi(0) \leq \mu \alpha \phi'(0), \qquad (2a)$$

that is,

$$f(\boldsymbol{x}_k + \alpha \boldsymbol{p}_k) - f(\boldsymbol{x}_k) \leq \mu \alpha \nabla f(\boldsymbol{x}_k) \boldsymbol{p}_k. \qquad (2b)$$

Figure 2: The interval ($R$) accepted by the Armijo step length rule.

# Typical convergence result

- *Suppose that $f \in C^1$, and that for the starting point $x_0$ it holds that the level set $\text{lev}_f(f(x_0)) = \{x \in \mathbb{R}^n \mid f(x) \leq f(x_0)\}$ is bounded. Consider the iterative algorithm defined on Page 1, with the following choices for each $k$:*
  - *$p_k$ satisfies the second sufficient descent condition on Page 7;*
  - *$\|p_k\| \leq M$, where $M$ is some positive constant; and*
  - *the Armijo step length rule is used.*

  *Then, the sequence $\{x_k\}$ is bounded, the sequence $\{f(x_k)\}$ is descending, lower bounded and therefore converges, and every limit point of $\{x_k\}$ is stationary.* □

- For convex $f$ every limit point is globally optimal.

# Step 4: Termination criteria

- Lesson number one: Cannot terminate based on the exact optimality conditions, because $\nabla f(\boldsymbol{x}) = \boldsymbol{0}^n$ exactly never happens!

- Unfortunate that we must compare with zero. Compare the lower bounding idea in Chapter 4.

- The recommendation is the combination of the following:
  1. $\|\nabla f(\boldsymbol{x}_k)\| \leq \varepsilon_1(1 + |f(\boldsymbol{x}_k)|)$, $\varepsilon_1 > 0$ small;
  2. $f(\boldsymbol{x}_{k-1}) - f(\boldsymbol{x}_k) \leq \varepsilon_2(1 + |f(\boldsymbol{x}_k)|)$, $\varepsilon_2 > 0$ small; and
  3. $\|\boldsymbol{x}_{k-1} - \boldsymbol{x}_k\| \leq \varepsilon_3(1 + \|\boldsymbol{x}_k\|)$, $\varepsilon_3 > 0$ small.

- Why? Need to cover cases of very steep and very flat functions.

- May need to use $\infty$-norm: $\|\boldsymbol{x}\|_\infty := \max_{1 \leq j \leq n} |x_j|$, for large $n$.

- Problem with the scaling of the problem: If

$$x_{k-1} = (1.44453, 0.00093, 0.0000079)^T,$$
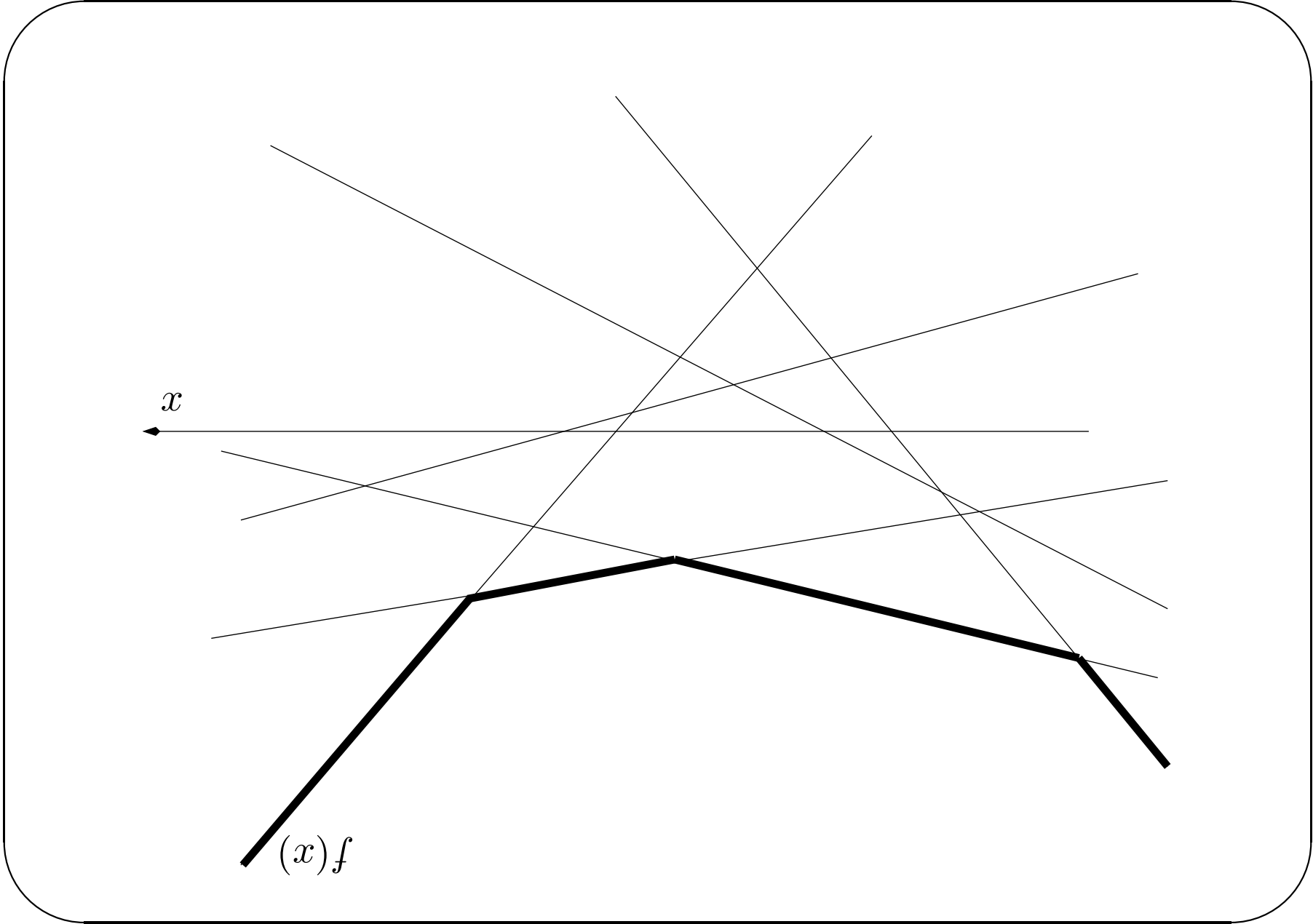
$$x_k = (1.44441, 0.00012, 0.0000011)^T;$$

$$\|x_{k-1} - x_k\|_\infty = \|(0.00012, 0.00081, 0.0000068)^T\|_\infty$$

$$= 0.00081.$$

- Small absolute error but large relative error!

- Better to apply the algorithm from a scaled problem where elements of $x$ have similar magnitude.

# Why is the $C^1$ property important?

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

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

- This is a piece-wise linear and convex function.

- It is differentiable almost everywhere, but *not* at the optimal solution!

- Ignoring non-differentiability may lead to the convergence to a non-optimal point.

- Convex functions always has *subgradients*, corresponding to all the possible slopes of the function.

- More on these when looking at Lagrangian duality!

# Trust region methods

- Trust region methods use quadratic models (as Newton).

- Avoids line searches by bounding the length of the search direction, at the same time influencing its direction.

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

- The model $\psi_k$ is *trusted* in a neighbourhood of $\boldsymbol{x}_k$ : $\|\boldsymbol{p}\| \leq \Delta_k$.

- Very useful when $\nabla^2 f(\boldsymbol{x}_k)$ is not positive semi-definite.

- Easy to minimize $\psi_k(\boldsymbol{p})$ subject to $\|\boldsymbol{p}\| \leq \Delta_k$.

- Idea: when $\nabla^2 f(\boldsymbol{x}_k)$ is badly conditioned, the value of $\Delta_k$ should be kept low (more of a steepest descent method); if $\nabla^2 f(\boldsymbol{x}_k)$ is well conditioned, $\Delta_k$ should become large to allow for unit steps (Newton! fast convergence).

- If $\Delta_k$ is small enough, $f(\boldsymbol{x}_k + \boldsymbol{p}_k) < f(\boldsymbol{x}_k)$ holds.

- Even if $\nabla f(\boldsymbol{x}_k) = \boldsymbol{0}^n$ holds, $f(\boldsymbol{x}_k + \boldsymbol{p}_k) < f(\boldsymbol{x}_k)$ still holds, if $\nabla^2 f(\boldsymbol{x}_k)$ is not positive definite.

- Progress from stationary points if saddle points or local maxima.

- Robust, quite popular.

- Update of trust region size based on a measure of similarity between the model $\psi_k$ and $f$: Let

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

If $\rho_k \leq \mu$ let $\quad \boldsymbol{x}_{k+1} = \boldsymbol{x}_k$ (unsuccessful step), else

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{p}_k \ \text{(successful step)}.$$

The value of $\mu$ is updated in the following manner, depending on the value of $p_k$:

$$p_k \leq \mu \implies \nabla_{k+1} = \frac{1}{2}\nabla_k,$$
$$\mu < p_k < n \implies \nabla_{k+1} = \nabla_k,$$
$$p_k \geq n \implies \nabla_{k+1} = 2\nabla_k.$$

# Minimizing implicit functions

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

$$x \in \mathbb{R}^n \longrightarrow \boxed{\text{Simulation}} \longrightarrow y \in \mathbb{R}^m$$

- 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 $x \mapsto f(x, y(x))$.

- Two distinct possibilities!

- (1) *Numerical differentiation* of $f$ by using a difference formula:

- Let $e_i = (0, 0, \ldots, 0, 1, 0, \ldots, 0)^T$ be the unit vector in $\mathbb{R}^n$. Then,

$$f(x + \alpha e_i) = f(x) + \alpha e_i^T \nabla f(x) + (\alpha^2/2) e_i^T \nabla^2 f(x) e_i + \cdots$$
$$= f(x) + \alpha \partial f(x)/\partial x_i + (\alpha^2/2) \partial^2 f(x)/\partial x_i^2 + \cdots$$

- So, for small $\alpha > 0$,

$$\frac{\partial f(x)}{\partial x_i} \approx \frac{f(x + \alpha e_i) - f(x)}{\alpha} \qquad \text{(forward difference)}$$

$$\frac{\partial f(x)}{\partial x_i} \approx \frac{f(x + \alpha e_i) - f(x - \alpha e_i)}{2\alpha} \qquad \text{(central difference)}$$

- Value of $\alpha$ typically set to a function of the machine precision; if too large, we get a bad approximation of the partial derivative, while a too small value might result in numerical cancellation.

- may work well if the simulation is accurate, otherwise bad derivative information.

- (2) *Derivative-free methods are available.* (Not counting subgradient methods, because they demand $f$ to be convex!) Either builds explicit *models* $\hat{f}$ of the objective function by evaluating $f$ at test points, or evaluates $f$ at grid points that are moved around, shrunk or expanded. *Names: Nelder–Mead, Pattern search.*