# Algorithms for optimization

Finding derivatives, unconstrained and constrained optimization

Zdeněk Hurák

Department of Control Engineering, Faculty of Electrical Engineering,
Czech Technical University in Prague

March 2, 2021

If possible, we want to use them (why shouldn't it be possible?).

How to compute them?
- symbolically
- numerically (finite difference, FD)
- algorithmic (also automatic) differentiation (AD)

Why?

$$\frac{\mathrm{d}f(x)}{\mathrm{d}x} \approx \frac{f(x + \alpha) - f(x)}{\alpha} \qquad \text{forward difference}$$

or

$$\frac{\mathrm{d}f(x)}{\mathrm{d}x} \approx \frac{f(x) - f(x - \alpha)}{\alpha} \qquad \text{backward difference}$$

or

$$\frac{\mathrm{d}f(x)}{\mathrm{d}x} \approx \frac{f(x + \frac{\alpha}{2}) - f(x - \frac{\alpha}{2})}{\alpha} \qquad \text{central difference}$$

Systematic application of chain rule for derivatives of composed functions.

How does it differ from symbolic and how from numerical differentiation?

Two versions: forward and reverse AD.

Similar to complex numbers, a dual number has two components:

$$x = v + d\epsilon, \quad \epsilon^2 = 0$$

Multiplication of two dual numbers $y = x_1 \cdot x_2$

$$x = (v_1 + d_1\epsilon) \cdot (v_2 + d_2\epsilon)$$
$$= v_1 v_2 + (v_1 d_2 + d_1 v_2)\epsilon$$

Matlab: recently added to Optimization Toolbox for Matlab,
CasADi, MAD (MatlabAD) from Tomlab, . . .

Python: CasADi, . . .

Julia: ForwardDiff.jl, ReverseDiff, Zygote.jl, . . .

- descent direction methods
- trust region methods

- finding a descent direction
- line search

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

Directional derivative negative

$$\nabla f(\mathsf{x}_k)^{\mathrm{T}}\mathsf{d}_k < 0$$

But beware of higher order terms – descent direction condition valid only in some vicinity of x.

$$d_k = -\nabla f(x_k)$$

$$\boxed{x_{k+1} = x_k - \alpha_k \nabla f(x_k)}$$

1. fixed step
2. exact search
3. approximate search

L-smoothness (Lipschitz continuity of the gradient)

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

If second derivative exists, $L$ is an upper bound

$$\|\nabla^2 f\| \leq L$$

For quadratic functions

$$L = \max \lambda_i(Q)$$

Quadratic dominance (descent lemma)

$$f(x_{k+1}) \leq f(x_k) + \nabla f(x_k)^\mathsf{T}(x_{k-1} - x_k) + \frac{L}{2}\|x_{k-1} - x_k\|^2$$

Step length

$$\boxed{\alpha = \frac{1}{L}}$$

Several methods (bisection, golden section, Newton, … )
For quadratic functions $f(\mathsf{x}) = \frac{1}{2}\mathsf{x}^\mathsf{T}\mathsf{Q}\mathsf{x} + \mathsf{c}^\mathsf{T}\mathsf{x}$ closed-form formula.

$$\underset{\alpha_k}{\text{minimize}}\, f(\mathsf{x}_k + \alpha_k \mathsf{d}_k)$$

$$
\begin{aligned}
f(\mathsf{x}_k + \alpha_k \mathsf{d}_k) &= \frac{1}{2}(\mathsf{x}_k + \alpha_k \mathsf{d}_k)^\mathsf{T}\mathsf{Q}(\mathsf{x}_k + \alpha_k \mathsf{d}_k) + \mathsf{c}^\mathsf{T}(\mathsf{x}_k + \alpha_k \mathsf{d}_k) \\
&= \frac{1}{2}\mathsf{x}_k^\mathsf{T}\mathsf{Q}\mathsf{x}_k + \mathsf{d}_k^\mathsf{T}\mathsf{Q}\mathsf{x}_k\alpha_k + \frac{1}{2}\mathsf{d}_k^\mathsf{T}\mathsf{Q}\mathsf{d}_k\alpha_k^2 + \mathsf{c}^\mathsf{T}(\mathsf{x}_k + \alpha_k \mathsf{d}_k)
\end{aligned}
$$

$$\frac{\mathrm{d}f(\mathsf{x}_k + \alpha_k \mathsf{d}_k)}{\mathrm{d}\alpha_k} = \mathsf{d}_k^\mathsf{T}(\mathsf{Q}\mathsf{x}_k + \mathsf{c}) + \mathsf{d}_k^\mathsf{T}\mathsf{Q}\mathsf{d}_k\alpha_k = 0$$

$$\boxed{\alpha_k = -\frac{\mathsf{d}_k^\mathsf{T}(\mathsf{Q}\mathsf{x}_k + \mathsf{c})}{\mathsf{d}_k^\mathsf{T}\mathsf{Q}\mathsf{d}_k} = -\frac{\mathsf{d}_k^\mathsf{T}\nabla f(\mathsf{x}_k)}{\mathsf{d}_k^\mathsf{T}\mathsf{Q}\mathsf{d}_k}}$$

Usually the exact minimum not needed, *sufficient descent* is enough.

Armijo / Wolfe conditions.

Backtracking algorithm: parameters $s$, $\beta \in (0,1)$, $\gamma \in (0,1)$:
Set $\alpha_k = s$
While

$$f(x_k) - f(x_k + \alpha_k d_k) < -\gamma \alpha_k d^T \nabla f(x_k),$$

set

$$\alpha_k = \beta \alpha_k.$$

$$\underset{x \in \mathbb{R}^2}{\text{minimize}} \, x^T Q x,$$

where the matrix $Q$ is

$$Q = \begin{bmatrix} 1000 & 20 \\ 20 & 1 \end{bmatrix}.$$

Condition number $\kappa$ for a given matrix A is

$$\kappa(A) = \|A^{-1}\| \cdot \|A\|.$$

It can be computed as ratio of the largest and smallest singular values, that is,

$$\kappa(A) = \frac{\sigma_{max}(A)}{\sigma_{min}(A)}.$$

Ideally should be around 1.
In the example above is well above 1000.

## Scaling to improve the conditioning and the convergence

Introduce new variable variable y

$$x = Sy,$$

The optimization cost changes $f(Sy)$. Relabel it to $g(y)$.
Chain rule

$$\nabla g(y) = S^T \nabla f(Sy).$$

Steepest descent iterations then change accordingly

$$y_{k+1} = y_k - \alpha_k \nabla g(y_k)$$

$$y_{k+1} = y_k - \alpha_k S^T \nabla f(Sy_k)$$

$$\underbrace{Sy_{k+1}}_{x_{k+1}} = \underbrace{Sy_k}_{x_k} - \alpha_k \underbrace{SS^T}_{D} \nabla f(\underbrace{Sy_k}_{x_k})$$

Defining the scaling matrix D as $SS^T$, a single iteration changes to

$$\boxed{x_{k+1} = x_k - \alpha_k D_k \nabla f(x_k).}$$

Make the Hessian matrix $\nabla^2 f(Sy)$ (the matrix Q above) better conditioned. Ideally, $\nabla^2 f(Sy) \approx I$.
Chain rule once again

$$\nabla^2 g(y) = S^T \nabla^2 f(Sy) S$$
$$= D^{\frac{1}{2}} \nabla^2 f D^{\frac{1}{2}}$$

A simple way using a diagonal scaling matrix D

$$D_{ii} = [\nabla^2 f(x_k)]_{ii}^{-1}.$$

Solve

$$g(x) = 0.$$

Approximate $g$ at $x_k$ using a linear function

$$\underbrace{g(x_{k+1})}_{0} = g(x_k) + g'(x_k)(x_{k+1} - x_k)$$

$$0 = g(x_k) + g'(x_k)x_{k+1} - g'(x_k)x_k,$$

from which the famous formula follows

$$\boxed{x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}.}$$

In the vector case

$$\boxed{x_{k+1} = x_k - J(x_k)^{-1}g(x_k).}$$

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

Model the function $f$ at $x_k$ using a quadratic function $m_k(x)$

$$m_k(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2} f''(x_k)(x - x_k)^2.$$

At the $k$-th iteration

$$\underset{x_{k+1} \in \mathbb{R}}{\text{minimize}} \quad m(x_{k+1})$$

Straightforward: find the value of $x_{k+1}$ for which the derivative of $m_k()$ vanishes.

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

Full vector version

$$\boxed{x_{k+1} = x_k - (\nabla^2 f(x_k))^{-1} \nabla f(x_k).}$$

$$x_{k+1} = x_k - \alpha_k \nabla^2 f(x_k)^{-1} \nabla f(x_k).$$

Generalization of the (scalar) secant method.
Secant approximation of the derivative (for rootfinding)

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

$$x_{k+1} = x_k - \underbrace{\frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})}}_{\approx \dot{f}(x_k)} f(x_k)$$

Secant approximation of the derivative (for optimization)

$$\ddot{f}(x_k) \approx \frac{\dot{f}(x_k) - \dot{f}(x_{k-1})}{x_k - x_{k-1}} =: b_k$$

$$b_k \underbrace{(x_k - x_{k-1})}_{s_{k-1}} = \underbrace{\dot{f}(x_k) - \dot{f}(x_{k-1})}_{y_{k-1}} \qquad \text{secant condition}$$

Secand condition in the vector case

$$\boxed{B_{k+1}s_k = y_k}$$

$B_k$ is a matrix with Hessian-like properties

$$B_k = B_k^\mathsf{T}$$

$$B_k \succ 0$$

How to get it? Updates.

$$B_{k+1} = B_k + \text{some "small" update}$$

Possibly updating $B_{k+1}^{-1}$ directly. One popular update is BFGS:

$$\boxed{H_{k+1} = H_k + \left(1 + \frac{y_k^\mathsf{T} H_k y_k}{s_k^\mathsf{T} y_k}\right) \cdot \frac{s_k s_k^\mathsf{T}}{s_k^\mathsf{T} y_k} - \frac{s_k y_k^\mathsf{T} H_k + H_k y_k s_k^\mathsf{T}}{y_k^\mathsf{T} s_k}}$$

Approximate $f()$ at $x_k$ with some model $m_k()$, typically a quadratic function

$$m_k(\mathsf{p}) = f(\mathsf{x}_k) + \nabla f(\mathsf{x}_k)^\mathsf{T} \mathsf{p} + \frac{1}{2}\mathsf{p}^\mathsf{T} \underbrace{\nabla^2 f(\mathsf{x}_k)}_{\text{or}\approx} \mathsf{p}$$

but trust the model only within

$$\|\mathsf{p}\|_2 \leq \delta$$

$$\begin{array}{|c|}
\hline
\underset{\mathsf{p}\in\mathbb{R}^n}{\text{minimize }} m_k(\mathsf{p}) \\
\text{subject to } \|\mathsf{p}\|_2 \leq \delta \\
\hline
\end{array}$$

and shrinking or expanding the trust region. Use

$$\eta = \frac{\text{actual improvement}}{\text{predicted improvement}} = \frac{f(x_k) - f(x_{k+1})}{f(x_k) - m_k(x_{k+1})}$$

Shrink for small $\eta$ ($\approx 0$) and expand for larger $\eta$ ($\approx 1$).