

# Quasi-Newton Methods in Optimisation

Andreas Langer

September 16, 2021

## Problem Description

Let  $n \in \mathbb{N}$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , the *objective function*.

Task: Find  $x^* \in \mathbb{R}^n$  such that  $f$  is minimal, i.e.,

$$f(x^*) = \min_{x \in \mathbb{R}^n} f(x)$$

Assumption:  $f$  is sufficiently smooth

# Characterisation of a Solution

A candidate  $x^* \in \mathbb{R}^n$  is called

(i) **global minimiser** of  $f$  in  $\mathbb{R}^n$ , if

$$f(x^*) \leq f(x) \quad \text{for all } x \in \mathbb{R}^n$$

(ii) **strict global minimiser** of  $f$  in  $\mathbb{R}^n$ , if

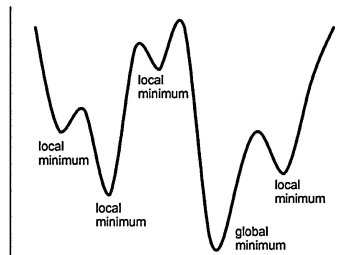
$$f(x^*) < f(x) \quad \text{for all } x \in \mathbb{R}^n \setminus \{x^*\}.$$

(iii) **local minimiser** of  $f$  in  $\mathbb{R}^n$ , if there exists a neighbourhood of  $x^*$  denoted by  $N(x^*)$  such that

$$f(x^*) \leq f(x) \quad \text{for all } x \in \mathbb{R}^n \cap N(x^*).$$

(iv) **strict local minimiser** of  $f$  in  $\mathbb{R}^n$ , if there exists  $N(x^*)$  such that

$$f(x^*) < f(x) \quad \text{for all } x \in (\mathbb{R}^n \cap N(x^*)) \setminus \{x^*\}.$$



## Existence of a Minimiser

In general a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  does not need to have a minimiser.

Reason:  $\mathbb{R}^n$  is **not** compact!

### Theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be continuous. If there exists an  $x_0 \in \mathbb{R}^n$  such that the level set

$$\mathcal{L}_f(x_0) := \{x \in \mathbb{R}^n \mid f(x) \leq f(x_0)\}$$

is compact, then there exists at least one global minimiser of  $f$  in  $\mathbb{R}^n$ .

# Notation

We denote:

- ▶ The gradient  $\nabla f(x) =: g(x) \in \mathbb{R}^n$  and write it as a row vector.
- ▶ The Hessian  $\nabla^2 f(x) =: G(x) \in \mathbb{R}^{n \times n}$  (it is a  $n \times n$  matrix).

Note:  $G$  is a symmetric matrix.

# Optimality Conditions

## Theorem (Second-order sufficient conditions)

A candidate  $x^* \in \mathbb{R}^n$  is a **local minimiser** of  $f$ , if

- ▶  $g(x^*) = 0$  (first-order necessary optimality condition)
- ▶ the Hessian  $G(x^*)$  is positive definite, i.e.,  $d^T G(x^*) d > 0$  for all  $d \in \mathbb{R}^n \setminus \{0\}$ .

A test for positive definiteness can be made together with Cholesky decomposition.  
See e.g. `scipy.linalg.chol`

If  $f$  is **convex**, then any **local minimiser** of  $f$  is also a **global minimiser** of  $f$ .

## Numerical Method: Newton

Solve  $g(x^*) = 0$  (first-order necessary optimality condition) by iterating:

- ▶ Choose an initial value (guess)  $x^{(0)} \in \mathbb{R}^n$
- ▶ Loop over  $k$  until a termination criterion holds:

$$\begin{aligned}s^{(k)} &:= -G(x^{(k)})^{-1}g(x^{(k)}) \\ x^{(k+1)} &:= x^{(k)} + s^{(k)}\end{aligned}$$

We write  $g^{(k)} := g(x^{(k)})$  and  $G^{(k)} := G(x^{(k)})$ .

$s^{(k)} := -G(x^{(k)})^{-1}g(x^{(k)})$  is called the *Newton direction*.

## Stopping Criterion

There are basically two termination criteria for Newton's method:

1. **Residual criterion:** The Newton iteration is stopped as soon as the residual  $\|g(x^{(k)})\|$  is small enough. In case of convergence we have

$$\lim_{k \rightarrow \infty} \|g(x^{(k)})\| = \|g(x^*)\| = 0.$$

2. **Cauchy criterion:** Terminate the iteration as soon as the Newton-correction  $\|x^{(k+1)} - x^{(k)}\| = \|s^{(k)}\|$  is small enough. In case of convergence we have

$$\lim_{k \rightarrow \infty} \|x^{(k+1)} - x^{(k)}\| = 0.$$

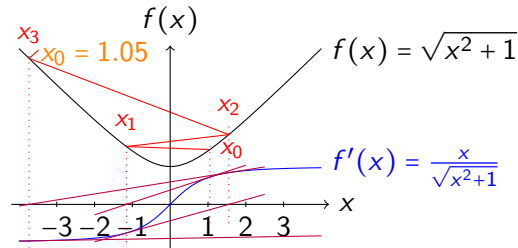
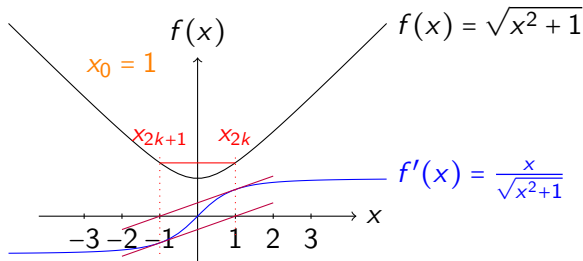
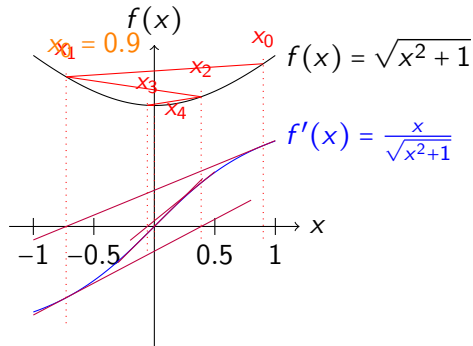
Both criteria should be *used with caution!*



# Newton Method: Problem (1)

1st Problem and its remedy:

- Local Convergence: Requires good initial guesses  $x^{(0)}$ .



Remedy: Globalization  $\rightarrow$  *Line search method*

## Newton Method: Problem (2)

2nd Problem and its remedy:

- ▶ Requires the evaluation of the Hessian  $G^{(k)}$  in each iteration.

Remedy: Choose a numerical approximation of  $G^{(k)}$  or even better of  $(G^{(k)})^{-1} \rightarrow$   
*Quasi Newton methods*

These two things lead to

$$x^{(k+1)} := x^{(k)} + \alpha^{(k)} s^{(k)}$$

with

$$s^{(k)} := -H^{(k)} g^{(k)}$$

where

- ▶  $\alpha^{(k)} > 0$  is a step size;
- ▶  $H^{(k)}$  is an approximation of  $(G^{(k)})^{-1}$ , which can be computed easily.

# Line Search

Determine  $\alpha^{(k)} > 0$  such that

$$f(x^{(k)} + \alpha^{(k)} s^{(k)}) < f(x^{(k)}).$$

- ▶ **Exact line search:**  $\alpha^{(k)} \in \arg \min_{\alpha \geq 0} f(x^{(k)} + \alpha s^{(k)})$
- ▶ **Inexact line search:** Armijo rule; Powell-Wolfe rule; Goldstein rule; ... (see e.g. [1],[2] in the course literature)

Define:  $\varphi(\alpha) := f(x^{(k)} + \alpha s^{(k)})$

We give two examples of such rules on the next slides.

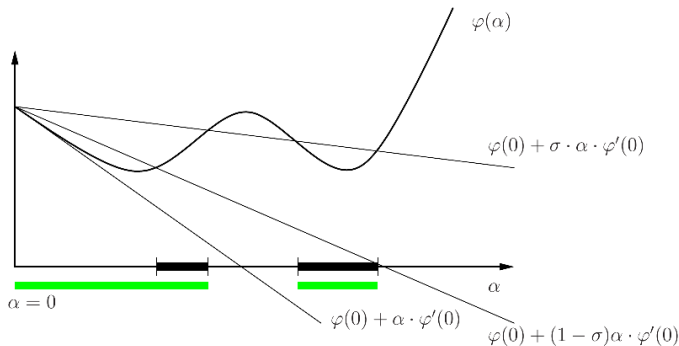
# Goldstein Rule

A step size  $\alpha$  is acceptable if the following two conditions

(i)  $\varphi(\alpha) \leq \varphi(0) + \sigma\alpha\varphi'(0)$  (Armijo rule)

(ii)  $\varphi(\alpha) \geq \varphi(0) + (1 - \sigma)\alpha\varphi'(0)$

hold for a given  $\sigma \in (0, \frac{1}{2})$ .



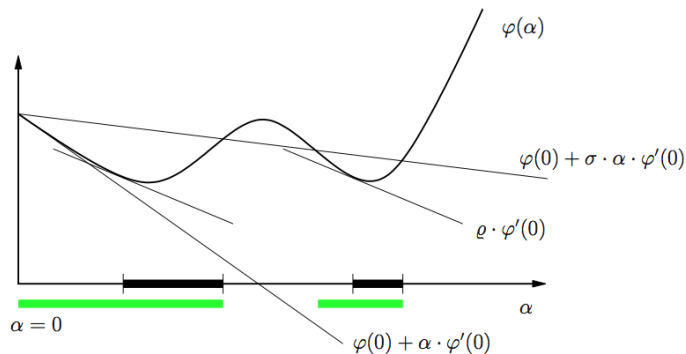
# Powell-Wolfe Rule

A step size  $\alpha$  is acceptable if the following two conditions

(i)  $\varphi(\alpha) \leq \varphi(0) + \sigma \alpha \varphi'(0)$  (Armijo rule)

(ii)  $\varphi'(\alpha) \geq \rho \varphi'(0)$

hold for given parameters  $\sigma \in (0, \frac{1}{2})$  and  $\rho \in (\sigma, 1)$ .



---

**Algorithm 1** Powell-Wolfe step size rule

---

Initialise: Parameters  $\sigma \in (0, \frac{1}{2})$ ,  $\rho \in (\sigma, 1)$ ,  $\alpha^- > 0$  (e.g.  $\sigma = 10^{-2}$ ,  $\eta = 0.9$ ,  $\alpha^- := 2$ )

**while**  $\alpha^-$  does not fulfil (Armijo) condition (i) **do**

$\alpha^- := \alpha^- / 2$

**end while**

Set  $\alpha^+ := \alpha^-$

**while**  $\alpha^+$  fulfils (Armijo) condition (i) **do**

$\alpha^+ := 2\alpha^+$

**end while**

**while**  $\alpha^-$  does not fulfil condition (ii) **do**

$\alpha_0 := \frac{\alpha^+ + \alpha^-}{2}$

**if**  $\alpha_0$  satisfies (Armijo) condition (i) **then**

$\alpha^- := \alpha_0$

**else**

$\alpha^+ := \alpha_0$

**end if**

**end while**

**return**  $\alpha := \alpha^-$

---

# Quasi-Newton Methods

A typical step looks like:

- ▶ Compute  $s^{(k)} := -H^{(k)}g^{(k)}$
- ▶ Perform line search to compute  $\alpha^{(k)}$ .
- ▶ Compute  $x^{(k+1)} := x^{(k)} + \alpha^{(k)}s^{(k)}$
- ▶ Update by *some* method  $H^{(k)} \rightarrow H^{(k+1)}$

## Motivation: Secant Method

In  $\mathbb{R}^1$ :

$$x^{(k+1)} := x^{(k)} - H^{(k)} g^{(k)}$$

with

$$G^{(k)} \approx \frac{g^{(k)} - g^{(k-1)}}{x^{(k)} - x^{(k-1)}} =: Q^{(k)} = H^{(k)-1}$$

In  $\mathbb{R}^1$  the approximation  $Q^{(k)}$  is **uniquely** determined by

$$Q^{(k)}(x^{(k)} - x^{(k-1)}) = g^{(k)} - g^{(k-1)}. \quad (1)$$

In  $\mathbb{R}^n$  is the Quasi-Newton condition (1) is **NOT uniquely** solvable.  
 $n$  equations for  $n^2$  unknowns  $Q_{ij}^{(k)}$

→ extra conditions needed.



# Broyden condition

Find  $Q^{(k)}$  by solving

$$\min \|Q^{(k)} - Q^{(k-1)}\|_F$$

subject to

$$Q^{(k)} \underbrace{(x^k - x^{k-1})}_{\delta^{(k)}} = \underbrace{g^{(k)} - g^{(k-1)}}_{\gamma^{(k)}}$$

This gives

$$Q^{(k)} = Q^{(k-1)} + \frac{\gamma^{(k)} - Q^{(k-1)}\delta^{(k)}}{\delta^{(k)\top}\delta^{(k)}}\delta^{(k)\top}$$

(see also "good Broyden's method" in [https://en.wikipedia.org/wiki/Broyden%27s\\_method](https://en.wikipedia.org/wiki/Broyden%27s_method))

## Sherman – Morrison (Woodbury) formula

Broyden update is a rank-1 update of the form

$$A_1 = A_0 + vw^T$$

where  $A_1, A_0 \in \mathbb{R}^{n \times n}$  and  $v, w \in \mathbb{R}^n$ .

Sherman – Morrison formula gives for the inverse

$$A_1^{-1} = A_0^{-1} - \frac{A_0^{-1}vw^TA_0^{-1}}{1 + w^TA_0^{-1}v}$$

if  $1 + w^TA_0^{-1}v \neq 0$ .

## Simple Rank-1 update

Consider the Sherman – Morrison formula above and replace  $A_0^{-1}$  by  $H^{(k-1)}$  and  $A_1^{-1}$  by  $H^{(k)}$ .

You then obtain:

$$H^{(k)} = H^{(k-1)} + \frac{(\delta^{(k)} - H^{(k-1)}\gamma^{(k)})}{\delta^{(k)\text{T}} H^{(k-1)} \gamma^{(k)}} \delta^{(k)\text{T}} H^{(k-1)}$$

# Broyden Condition for Inverse Hessian

Alternatively, we might approximate the inverse Jacobian directly by

Find  $H^{(k)} := Q^{(k)-1}$  by solving

$$\min \|H^{(k)} - H^{(k-1)}\|_F$$

subject to

$$Q^{(k)} \underbrace{(x^k - x^{k-1})}_{\delta^{(k)}} = \underbrace{g^{(k)} - g^{(k-1)}}_{\gamma^{(k)}}$$

This gives

$$H^{(k)} = H^{(k-1)} + \frac{\delta^{(k)} - H^{(k-1)}\gamma^{(k)}}{\gamma^{(k)\text{T}}\gamma^{(k)}}\gamma^{(k)\text{T}}$$

(see also "bad Broyden's method" in [https://en.wikipedia.org/wiki/Broyden%27s\\_method](https://en.wikipedia.org/wiki/Broyden%27s_method))

## Symmetric Rank 1

Start with a symmetric and invertible matrix  $Q^{(k)}$ .

**Task:** Find a symmetric  $Q^{(k+1)}$  via a rank-1 update, i.e., of the form

$$Q^{(k+1)} = Q^{(k)} + vw^T,$$

fulfilling the Quasi-Newton condition.

This gives

$$Q^{(k)} = Q^{(k-1)} + \frac{(\gamma^{(k)} - Q^{(k-1)}\delta^{(k)})(\gamma^{(k)} - Q^{(k-1)}\delta^{(k)})^T}{(\gamma^{(k)} - Q^{(k-1)}\delta^{(k)})^T \delta^{(k)}}$$

## Symmetric Rank 1 (Inverse)

Consider the Sherman – Morrison formula above and replace  $A_0^{-1}$  by  $H^{(k-1)}$  and  $A_1^{-1}$  by  $H^{(k)}$ .

You then obtain:

$$H^{(k)} = H^{(k-1)} + a u u^T$$

with

$$u := \delta^{(k)} - H^{(k-1)} \gamma^{(k)}, \quad a := \frac{1}{u^T \gamma^{(k)}}$$

## Rank-2 Update – DFP Method

Davidson-Fletcher-Powell (DFP) update

$$Q^{(k+1)} := Q^{(k)} + \left(1 + \frac{\delta^{(k)\top} Q^{(k)} \delta^{(k)}}{\gamma^{(k)\top} \delta^{(k)}}\right) \frac{\gamma^{(k)} \gamma^{(k)\top}}{\gamma^{(k)\top} \delta^{(k)}} - \frac{\gamma^{(k)} \delta^{(k)\top} Q^{(k)} + Q^{(k)} \delta^{(k)} \gamma^{(k)\top}}{\gamma^{(k)\top} \delta^{(k)}}$$

$$H^{(k+1)} := H^{(k)} + \frac{\delta^{(k)} \delta^{(k)\top}}{\delta^{(k)\top} \gamma^{(k)}} - \frac{H^{(k)} \gamma^{(k)} \gamma^{(k)\top} H^{(k)}}{\gamma^{(k)\top} H^{(k)} \gamma^{(k)}}$$

## Rank-2 update - BFGS method

Broyden-Fletcher-Goldfarb-Shanno (BFGS) update

$$Q^{(k+1)} := Q^{(k)} + \frac{\gamma^{(k)}\gamma^{(k)\text{T}}}{\gamma^{(k)\text{T}}\delta^{(k)}} - \frac{Q^{(k)}\delta^{(k)}\delta^{(k)\text{T}}Q^{(k)}}{\delta^{(k)\text{T}}Q^{(k)}\delta^{(k)}}$$

$$H^{(k+1)} := H^{(k)} + \left(1 + \frac{\gamma^{(k)\text{T}}H^{(k)}\gamma^{(k)}}{\delta^{(k)\text{T}}\gamma^{(k)}}\right) \frac{\delta^{(k)}\delta^{(k)\text{T}}}{\delta^{(k)\text{T}}\gamma^{(k)}} - \frac{\delta^{(k)}\gamma^{(k)\text{T}}H^{(k)} + H^{(k)}\gamma^{(k)}\delta^{(k)\text{T}}}{\delta^{(k)\text{T}}\gamma^{(k)}}$$

see, e.g., Fletcher, R: Practical Methods of Optimization, 2nd Ed, p.55 (reference [1] in the course literature)