

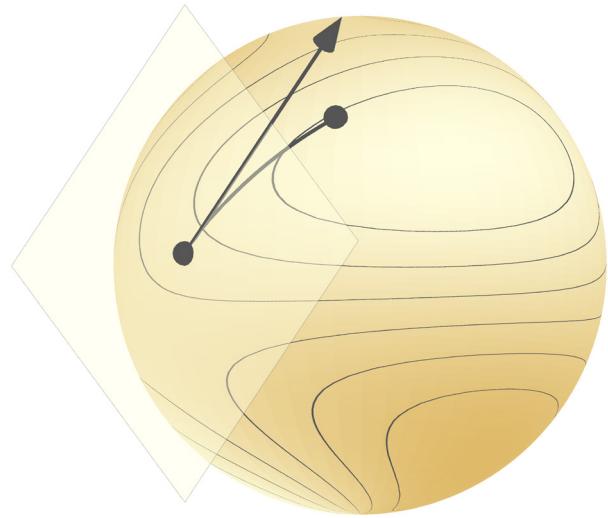
208

# Computing the Newton step

Spring 2023

Optimization on manifolds, MATH 512 @ EPFL

Instructor: Nicolas Boumal



# The linear equation of Newton steps

Let  $f: \mathcal{M} \rightarrow \mathbf{R}$  be smooth on a Riemannian manifold  $\mathcal{M}$ .

$$Ax = b$$

A:  $n \times n$

$n \times 1$

The Newton step at  $x$  is the tangent vector  $s \in T_x \mathcal{M}$  which solves:

$$\text{Hess } f(x) : T_x \mathcal{M} \rightarrow T_x \mathcal{M}$$

$$\text{Hess } f(x)[s] = -\text{grad } f(x)$$

$$\text{Hess } f(x) \succ 0$$

This is a **linear** equation in the unknown  $s$ . How do we solve it?

(Once solved, move from  $x$  to  $R_x(s)$  with the chosen retraction  $R$ .)

$$H: T_x M \rightarrow T_x M : H = \text{Hess}f(x)$$

$$b \in T_x M \quad , \quad b = -\text{grad}f(x)$$

# Matrix-based approaches

Pick a basis  $u_1, \dots, u_d \in T_x M$  for  $T_x M$ .

We have a linear system to solve:  $HA = b$ , for  $a \in T_x M$ .

$$A = \sum_{i=1}^d \alpha_i u_i ; \quad \sum_i \alpha_i Hu_i = b$$

$\overline{H}_{ij} = \langle Hu_i, u_j \rangle_x$ $b_j = \langle b, u_j \rangle_x$ $F_A = b$	$\left  \begin{array}{l} \left\langle \sum_i \alpha_i Hu_i, u_j \right\rangle_x = \langle b, u_j \rangle_x \quad \forall j \\ \sum_i \alpha_i \langle Hu_i, u_j \rangle_x = \langle b, u_j \rangle_x \quad \forall j \end{array} \right.$
--	---

Paradigm: We have access to  $\text{Hess}f(x)$  as an **operator**.

Calling that operator dominates computational costs.

# Matrix-free approaches

Assuming  $\text{Hess}f(x) > 0$ , recall the Newton step  $s$  **minimizes**:

$$m_x(v) = \frac{1}{2} \langle v, \text{Hess}f(x)[v] \rangle_x + \langle \text{grad}f(x), v \rangle_x + f(x)$$

Maybe we could run an optimization algorithm on  $m_x: T_x \mathcal{M} \rightarrow \mathbb{R}$ ?

$$g: T_x \mathcal{M} \rightarrow \mathbb{R}, \quad g(v) = \frac{1}{2} \langle v, H_v \rangle_x - \langle b, v \rangle_x, \quad \begin{cases} b = -\text{grad}f(x) \\ H = \text{Hess}f(x) \end{cases}$$

Gradient descent to  $g$ ?

$$v_0 = 0; \quad v_n = v_{n-1} - d_n \text{grad}g(v_{n-1})$$

$$\underline{\text{grad } g(v) = Hv - b}$$

$$g(v + \alpha r) = \frac{1}{2} \langle v + \alpha r, Hv + \alpha r \rangle_x - \langle b, v + \alpha r \rangle_x$$

$$= \frac{1}{2} \langle v, Hv \rangle_x + \frac{\alpha}{2} \langle v, Hr \rangle_x + \frac{\alpha}{2} \langle r, Hv \rangle_x$$

$$+ \frac{\alpha^2}{2} \langle r, Hr \rangle_x - \langle b, v \rangle_x - \alpha \langle b, r \rangle_x$$

$$= g(v) + \alpha \langle r, Hv - b \rangle_x + \frac{\alpha^2}{2} \langle r, Hr \rangle_x$$

$$\Rightarrow \text{optimal } \alpha \text{ is: } \alpha = -\frac{\langle r, Hv - b \rangle_x}{\langle r, Hr \rangle_x}$$

$$\underline{\text{For GD: } \alpha = \frac{\|r\|_x^2}{\langle r, Hr \rangle_x}}$$

Gradient descent to minimize  $g(v) = \frac{1}{2} \langle v, Hv \rangle_x - \langle b, v \rangle_x$  on  $T_x \mathcal{M}$ :

Initialize  $v_0 = 0$  and  $r_0 = -\text{grad } g(v_0) = b - Hv_0 = b$

For  $n$  in 1, 2, 3, ...

- Compute  $Hr_{n-1}$ . (only call to  $H$  per iteration)
- $\alpha_n = \frac{\|r_{n-1}\|_x^2}{\langle r_{n-1}, Hr_{n-1} \rangle_x}$
- $v_n = v_{n-1} + \alpha_n r_{n-1}$
- $r_n = -\text{grad } g(v_n) = b - Hv_n = b - H(v_{n-1} + \alpha_n r_{n-1}) = b - Hv_{n-1} - \alpha_n Hr_{n-1} = r_{n-1} - \alpha_n Hr_{n-1}$
- If  $\|r_n\|_x \leq \text{tol} \cdot \|b\|_x$ , output  $v_n$

GD is ok, but we can surely do better.

Observe that GD iterates are linear combinations of the residues:

$$\begin{aligned} v_n &= v_{n-1} + \alpha_n r_{n-1} = v_{n-2} + \alpha_{n-1} r_{n-2} + \alpha_n r_{n-1}, \\ &\vdots \\ &= \alpha_1 r_0 + \alpha_2 r_1 + \dots + \alpha_n r_{n-1}. \end{aligned}$$

Idea: Choose  $\alpha_1, \dots, \alpha_n$  s.t.  $v_n$  minimizes  $g(v) = \frac{1}{2} \langle v, Hv \rangle_a - \langle b, v \rangle_a$   
over the whole subspace spanned by  $r_0, \dots, r_{n-1}$ .

$$\begin{aligned} g\left(\sum_{i=1}^n \alpha_i r_{i-1}\right) &= \frac{1}{2} \left\langle \sum_i \alpha_i r_{i-1}, H\left(\sum_j \alpha_j r_{j-1}\right) \right\rangle_a - \left\langle b, \sum_i \alpha_i r_{i-1} \right\rangle_a \\ &= \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \underbrace{\left\langle r_{i-1}, Hr_{j-1} \right\rangle_a}_{-\sum_i \alpha_i \left\langle b, r_{i-1} \right\rangle_a} \end{aligned}$$

$$= \frac{1}{2} \alpha^T M \alpha - c^T \alpha.$$

$M > 0$  as long as the vectors  $r$  are linearly independent.

$\Rightarrow$  optimal  $\alpha$  solves  $M\alpha = c$ .

optimal :  $v_n = \sum_{i=1}^n \alpha_i p_{i-1}$

Gram-Schmidt theorem

Vectors  $r_0, r_1, \dots$

with:  $M\alpha = c$ ,

$$\Pi_{ij} = \langle p_{i-1}, H p_{j-1} \rangle_x$$

$$\alpha_i = \frac{\langle b, p_{i-1} \rangle_x}{\langle p_{i-1}, H p_{i-1} \rangle_x}$$

## Gradient descent (GD)

Initialize  $v_0 = 0$ ,  $r_0 = b$

For  $n$  in 1, 2, 3, ...

- Compute  $Hr_{n-1}$
- $\alpha_n = \frac{\|r_{n-1}\|_x^2}{\langle r_{n-1}, Hr_{n-1} \rangle_x}$
- $v_n = v_{n-1} + \alpha_n r_{n-1}$
- $r_n = r_{n-1} - \alpha_n Hr_{n-1}$
- If  $\|r_n\|_x \leq \text{tol} \cdot \|b\|_x$ , output  $v_n$

$$r_n = -\text{grad } g(v_n) = b - Hv_n$$

## Conjugate gradients (CG)

Initialize  $v_0 = 0$ ,  $r_0 = b$ ,  $p_0 = r_0$

For  $n$  in 1, 2, 3, ...

- Compute  $Hp_{n-1}$
  - $\alpha_n = \frac{\|r_{n-1}\|_x^2}{\langle p_{n-1}, Hp_{n-1} \rangle_x}$
  - $v_n = v_{n-1} + \alpha_n p_{n-1}$
  - $r_n = r_{n-1} - \alpha_n Hp_{n-1}$
  - If  $\|r_n\|_x \leq \text{tol} \cdot \|b\|_x$ , output  $v_n$
  - $\beta_n = \frac{\|r_n\|_x^2}{\|r_{n-1}\|_x^2}$
  - $p_n = r_n + \beta_n p_{n-1}$
- Cheap Gram-Schmidt*

$$p_n = r_n - \sum_{i=0}^{n-1} \frac{\langle r_n, H p_i \rangle}{\langle p_i, H p_i \rangle} p_i$$

$$g(\nu) = \frac{1}{2} \langle \nu, H\nu \rangle_x - \langle b, \nu \rangle_x$$

# Finite termination of CG...

CG iterates are  $\nu_n = \underset{\nu \in \mathcal{K}_n}{\operatorname{argmin}} g(\nu)$  with  $\mathcal{K}_n = \operatorname{span}(r_0, \dots, r_{n-1})$ .

Assume  $r_0, \dots, r_{n-1}$  are linearly independent (for some  $n$ ).

If  $r_n = 0$ , then  $\nu_n$  is the global optimum: stop.

If  $r_n \neq 0$ , then  $r_0, \dots, r_n$  are linearly independent.

Since  $\mathcal{K}_n \subseteq T_x \mathcal{M}$ , unless we stop early,  $\mathcal{K}_n = T_x \mathcal{M}$  for  $n = \dim \mathcal{M}$ .

Thus,  $\nu_n$  is the exact solution for some  $n \leq \dim \mathcal{M}$ .

... is a mathematical fairy tale.

Numerically,  $p_0, p_1, p_2, \dots$  progressively lose  $H$ -orthogonality.

As a result, implementations of CG do *not* find the exact solution.

But CG is still much better than GD!

In theory *and* in practice, both GD and CG enjoy linear convergence.

Only, CG's rate is typically much better than GD's.

Enough to get machine precision (or stop early).