

Conjugate gradient methods = cg – Methods

1. Conjugate directions -definition
2. Orthogonalization
3. cg-method for strictly convex quadratic functions
4. Preconditioning
5. cg-method for smooth non quadratic functions
6. Application
7. Experiments

Orthogonalization by E. Schmidt (recall)

Algorithm A7 (GRAM-SCHMIDT)

S0. Given linear independent $\mathbf{p}^0, \dots, \mathbf{p}^{m-1} \in \mathbb{R}^n, m \leq n-1,$
 $\mathbf{Q} \in \text{SPD}^n,$ set $\mathbf{d}^0 := \mathbf{p}^0, k := 0$

S1. while $k < m$

$$\mathbf{d}^k = \mathbf{p}^k - \sum_{i=0}^{k-1} \beta_{ik} \mathbf{d}^i \quad (\mathbf{Q} \text{ orthogonal projection})$$

with

$$\beta_{ik} := \frac{(\mathbf{p}^k)^T \mathbf{Q} \mathbf{d}^i}{(\mathbf{d}^i)^T \mathbf{Q} \mathbf{d}^i} \quad (\text{Fourier's coefficients})$$

S2. Set $k := k + 1$ and goto S1.

Conjugate directions -definition

Bijjective affine mapping:

circle \cup square of tangents \cup cross of diameters

orthogonal to each other

orthogonal w.r.t. usual scalar product

$$\mathbf{x} \bullet \mathbf{y} = \mathbf{x}^T \mathbf{E} \mathbf{y}$$

↓

Ellipsis \cup parallelogram of tangents \cup cross of diameters

conjugate to each other

orthogonal w.r.t. a new skalarproduct

$$\mathbf{x} \bullet \mathbf{y} = \mathbf{x}^T \mathbf{Q} \mathbf{y}, \mathbf{Q} \in \text{SPD}$$

Def.: $\mathbf{Q} \in \text{SPD},$

\mathbf{x} \mathbf{Q} -conjugate (\mathbf{Q} -orthogonal) to $\mathbf{y} : \iff \mathbf{x}^T \mathbf{Q} \mathbf{y} = 0$

Minimal property of variety $\mathbf{x}^0 + \mathcal{L}\{\mathbf{d}^0, \dots, \mathbf{d}^k\}$

Theorem:cf. [Reinhardt et.al. Th. 3.68] Assume:

1. $f(\mathbf{x}) := \mathbf{b}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x}, \mathbf{Q} \in \text{SPD}^n.$
2. $\mathbf{d}^0, \dots, \mathbf{d}^{n-1}$ \mathbf{Q} -conj. system and $\mathbf{x}^0 \in \mathbb{R}^n$
3. $\mathbf{x}^{k+1} = \mathbf{x}^k + \lambda_k \mathbf{d}^k$ for $k = 0, \dots, n-1$
 $\lambda_k > 0$ perfect step length for $f(\mathbf{x}^k + \lambda \mathbf{d}^k) \rightarrow \min$
4. $V_{k+1} := \{\mathbf{x} \mid \mathbf{x} = \mathbf{x}^0 + \sum_{i=0}^k \mu_i \mathbf{d}^i, \mu_i \in \mathbb{R}\}$

Then hold (!! $\nabla f(\mathbf{x}^{k+1}) = \mathbf{Q} \mathbf{x}^{k+1} + \mathbf{b} = \mathbf{Q} \mathbf{x}^k + \mathbf{b} + \lambda \mathbf{Q} \mathbf{d}^k$)

- i) $\nabla f(\mathbf{x}^{k+1})^T \mathbf{d}^j = 0$ for $0 \leq j \leq k, k = 0, 1, \dots, n-1$
- ii) $\mathbf{x}^{k+1} = \text{argmin}\{f(\mathbf{x}) \mid \mathbf{x} \in V_{k+1}\}, k = 0, 1, \dots, n-1$
- iii) $\nabla f(\mathbf{x}^{i+1}) - \nabla f(\mathbf{x}^i) = \lambda_i \mathbf{Q} \mathbf{d}^i, i = 0, 1, \dots, n-1$

Algorithm A8 (cg-method for strictly convex quadratic functions)

- S0. $\mathbf{x}^0 \in \mathbb{R}^n$, $\mathbf{d}^0 := -\nabla f(\mathbf{x}^0)$, $\varepsilon > 0$, $k := 0$
- S1. **if** $\|\nabla f(\mathbf{x}^k)\| = 0$ ($< \varepsilon$ on computer), **STOPP**
- S2. **perfect step length:** $\lambda_k := -\frac{\nabla f(\mathbf{x}^k)^T \mathbf{d}^k}{(\mathbf{d}^k)^T \mathbf{Q} \mathbf{d}^k}$
 update: $\mathbf{x}^{k+1} := \mathbf{x}^k + \lambda_k \mathbf{d}^k$
- S3. **New descent direction: (!! only last summand)**
 $\mathbf{d}^{k+1} := -\nabla f(\mathbf{x}^{k+1}) + \beta_k \mathbf{d}^k$ with $\beta_k = \frac{\nabla f(\mathbf{x}^{k+1})^T \mathbf{Q} \mathbf{d}^k}{(\mathbf{d}^k)^T \mathbf{Q} \mathbf{d}^k}$
- S4. Set $k := k + 1$ and goto S1.

Numerical realization of A8

- S0. $\mathbf{x}^0 \in \mathbb{R}^n$, $\mathbf{g}^0 := \nabla f(\mathbf{x}^0)$, $\mathbf{d}^0 = -\mathbf{g}^0$, $\varepsilon > 0$, $k := 0$
- S1. **if** $\|\mathbf{g}^k\| = 0$ ($< \varepsilon$ on computer), **STOPP**
- S2. **perfect step length:** $\lambda_k := \frac{\|\mathbf{g}^k\|^2}{(\mathbf{d}^k)^T \mathbf{Q} \mathbf{d}^k}$
 $\mathbf{x}^{k+1} := \mathbf{x}^k + \lambda_k \mathbf{d}^k$
 $\mathbf{g}^{k+1} := \mathbf{g}^k + \lambda_k \mathbf{Q} \mathbf{d}^k (= \nabla f(\mathbf{x}^{k+1}))$
- S3. **New descent direction:**
 $\mathbf{d}^{k+1} := -\mathbf{g}^{k+1} + \beta_k \mathbf{d}^k$ with $\beta_k = \frac{\|\mathbf{g}^{k+1}\|^2}{\|\mathbf{g}^k\|^2}$
- S4. Set $k := k + 1$ and goto S1.

Calculate \mathbf{g}^{k+1} sometimes directly (error transmission!!)

Theorem:cf. [Reinhardt et al. Th. 3.69] Let \mathbf{x}^k , $k = 0, 1, \dots, m + 1 \leq n$ generated by Algorithm A8. Then (a)-(f) hold for $k = 1, 2, \dots, m$

- (a) $\nabla f(\mathbf{x}^j)^T \nabla f(\mathbf{x}^i) = 0$, $0 \leq i < j \leq k$ OGS
- (b) $(\mathbf{d}^j)^T \mathbf{Q} \mathbf{d}^i = 0$, $0 \leq i < j \leq k$, **Q**-OGS
- (c) $\nabla f(\mathbf{x}^j)^T \mathbf{d}^i = 0$, $0 \leq i < j \leq k$ minimality w.r.t. V_k
- (d) $\nabla f(\mathbf{x}^j)^T \mathbf{d}^j = -\|\nabla f(\mathbf{x}^j)\|^2$, $0 \leq j \leq k$, **descent cond.(DC)**
- (e) $V_k - \mathbf{x}^0 = \mathcal{L}\{\mathbf{d}^0, \dots, \mathbf{d}^{k-1}\} = \mathcal{L}\{\nabla f(\mathbf{x}^0), \dots, \nabla f(\mathbf{x}^{k-1})\}$
- (f) $V_k - \mathbf{x}^0 = \mathcal{L}\{\nabla f(\mathbf{x}^0), \mathbf{Q} \nabla f(\mathbf{x}^0), \dots, \mathbf{Q}^{k-1} \nabla f(\mathbf{x}^0)\}$
- (g) $\nabla f(\mathbf{x}^j)^T \mathbf{Q} \mathbf{d}^i = 0$, $0 \leq i < j - 1 \leq k - 1$ for $k = 2, 3, \dots, m$

The algorithm A8 stops after $k \leq n$ steps at $\mathbf{x}^k = \mathbf{x}^*$ (**finite stop property**)

Lemma: BFGS / DFP - method, f quadratic, strictly convex, **perfect step length** \implies

- \mathbf{d}^k create a **Q**-OGS ($\implies \mathbf{x}^* = \mathbf{x}^k$, $k \leq n$)
- $\mathbf{H}_0 = \mathbf{E} \implies$ BFGS (\equiv DFP) \equiv CG

Structure of eigenvalues

Theorem: [Bertsekas(1999), Pytlak(2009), Reinhardt et. al. Th. 3.72] Given f quadratic, strictly convex with minimum point \mathbf{x}^* and $\text{EV}(\mathbf{Q})$: $\lambda_1, \dots, \lambda_n \in \mathbb{R}$, iterations \mathbf{x}^k , $k = 0, 1, \dots, m$ w.r.t. Algorithm A8 such that $\mathbf{x}^m = \mathbf{x}^*$. Then the estimation

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \leq \max_{j=1, \dots, n} (P_{k+1}(\lambda_j))^2 (f(\mathbf{x}^0) - f(\mathbf{x}^*))$$

is valid for $k = 0, \dots, m - 1$ and **each** Polynom $P_{k+1}(\lambda)$ of degree $k + 1$ with $P_{k+1}(0) = 1$.

Corollary: If $0 < a < \lambda_1 \leq \dots \leq \lambda_{n-k} \leq b < \lambda_{n-k+1} \leq \dots \leq \lambda_n$ then

$$f(\mathbf{x}^{k+1}) - f(\mathbf{x}^*) \leq \left(\frac{b-a}{b+a}\right)^2 (f(\mathbf{x}^0) - f(\mathbf{x}^*))$$

Corollary: If only $r \leq n$ of eigenvalues of \mathbf{Q} are pairwise disjoint then

$$\mathbf{x}^r = \mathbf{x}^*$$

Preconditioning – intention: $B \approx Q$

S regular, Transformation: $Sy := x, f_S(y) := f(Sy)$,
 chain rule: $\nabla f_S(y) = S^T \nabla f(x)$. cg-method for f_S yields with
 $B^{-1} := SS^T$ the cg-method with preconditioning in x -data

Algorithm A8 (cg-method for strictly convex quadratic functions)

S0. $x^0 \in \mathbb{R}^n, g^0 := \nabla f(x^0), d^0 = -B^{-1}g^0, \varepsilon > 0, k := 0$

S1. if $\|g^k\| = 0$ ($< \varepsilon$ on computer), STOPP

S2. transformed perfect step length: $\lambda_k := \frac{(g^k)^T B^{-1}g^k}{(d^k)^T Qd^k}$

$x^{k+1} := x^k + \lambda_k d^k$ (instead $B^{-1}g := h$ solve fast (!) $Bh = g$)

$g^{k+1} := g^k + \lambda_k Qd^k (= \nabla f(x^{k+1}))$

S3. New descent direction:

$d^{k+1} := -B^{-1}g^{k+1} + \beta_k d^k$ with $\beta_k = \frac{(g^{k+1})^T B^{-1}g^{k+1}}{(g^k)^T B^{-1}g^k}$

S4. Set $k := k + 1$ and goto S1.

NLO-adaption

Use of identical reformulation in the quadratic case

$$\beta_k = \frac{\nabla f(x^{k+1})^T Qd^k}{(d^k)^T Qd^k} \quad f \text{ quadratic}$$

smooth NLO-problem

$$\stackrel{\text{i)+d) + iii)}}{=} \frac{\nabla f(x^{k+1}) (\nabla f(x^{k+1}) - \nabla f(x^k))}{\|\nabla f(x^k)\|^2} \quad \text{POLAK/RIBIERE}$$

$$= \frac{\|\nabla f(x^{k+1})\|^2}{\|\nabla f(x^k)\|^2} \quad \text{FLETCHER/REEVES}$$

- Perf. step length. \mapsto AR with IP., PW, Backtracking,...
the more perfect the better, but effort observe
- There are: **the** cg for QP, but **a lot of** cg for NLO

Examples for Preconditioning with $B \approx Q$

1. incomplete Cholesky decomposition preconditioner (several strategies): $B = L^T L$

Application: $L^T y = g, Lh = y \Leftrightarrow h = B^{-1}g$

fminunc: Chol. decomp. for r -diagonals of Q

EdOptLab: cholinc with drop param. $p \in (0, \infty)$

2. SSOR (symmetric successive overrelaxation, $\omega \in [1, 2)$) preconditioner:

$B = (D + \omega L)D^{-1}(D + \omega L)^T / (2 - \omega)$ where

$L_{ij} = Q_{ij}, i > j; L_{ij} = 0$ elsewhere, $D = \text{diag}(Q)$

Application: $(D + \omega L)u = g; w = Du;$

$(D + \omega L)^T h = (2 - \omega)w$

3. Jacobi-preconditioner: $B = \text{diag}(Q)$

Application: $h_i = (Q_{ii})^{-1}g_i$

About convergence theory

- Poljak/Ribiere: $f \in C^2, f''$ unif. SPD, λ_k perfect \implies R-linearly convergent [SCHW79, S. 236]
- Fletcher/Reeves: under same assumptions only convergent [Geig90, S. 231]
- F/R: $c \leq f \in C^{1,L}, PW \implies$ culm. point (CP) is stationary [Geig90, S. 229]
- P/R $c \leq f \in C^{1,L}, PW, \|x^{k+1} - x^k\| \rightarrow 0 \implies$ CP is stationary [Geig90, S. 232]
- P/R, F/R: $f \in C^3, f''$ unif. SPD, λ_k perfect $\implies n$ -step quadratic. conv. with restart [COHEN 72], [Schw79, S. 238]
- F/R: $d^k \rightarrow$ parallel \implies restart with $d^k = -\nabla f(x^k)$

Application of cg-methods

- approximate solution of linear Newton-equation, see inexact Newton methods and TR-methods
- large scale NLO using preconditioning at periodical restart [Pytlak 2009], faster than BFGS and limited memory BFGS, sophisticated algorithms
- (large scale) strictly convex QP with preconditioning

Experiments with F/R, F/R Restart and P/R

- CG14: non perfect step length is bad
- CG03: Comparison CG-PR, CG-FR, CG-Q, BFGS, Newton for strictly convex quadratic function till $\dim x = 100$ for $\kappa(Q) = 100$
- CG04: As CG03 but with $\kappa(Q) = 10000$
- CG06, CG07, CG09, CG10: BFGS for strictly convex quadratic function with different preconditioning ($H_0 = B$)
- CG11: CG-Q with preconditioning
- CG12: Comparison CG-PR, CG-FR, BFGS, Banana function, perfect step length
- CG15: CG-methods with several restarts
- CG16: Effectivity of CG-PR, CG-FR(restart), CG-PR(restart), BFGS for Banana till dimension 100 with parameter $\alpha = 100$ and $\alpha = 1$ (simulation of preconditioning)