

Conjugate Gradient method

Sunday, March 21, 2021

8:49 AM

Fast algo. for smooth unconstrained optim.

- grad. descent, Nesterov
- Newton (see Boyd + Vandenberghe, self-concordant analysis)
- Nonlinear conjugate gradient

* - Quasi-Newton (BFGS, L-BFGS)

• (linear) Conjugate Gradient (these notes follow Nocedal + Wright text and Anne Greenbaum's monograph)

History: Hestenes, Stiefel '50s (nonlinear: Fletcher + Reeves '60s)

CG solves $Ax=b$ if $A \succ 0$, A is $n \times n$
approximately (naively: $O(n^3)$ flops)

Motivation / link to optimization:

$$\min_x \frac{1}{2} \|\tilde{A}x - \tilde{b}\|^2$$

$$\varphi(x) = \frac{1}{2} x^T \underbrace{\tilde{A}^T \tilde{A}}_A x - \underbrace{\tilde{b}^T \tilde{A}}_{b^T} x + \frac{1}{2} \tilde{b}^T \tilde{b}$$

$$= \frac{1}{2} x^T A x - b^T x + \text{const.}$$

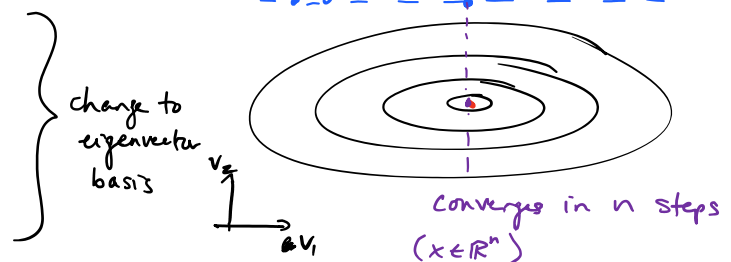
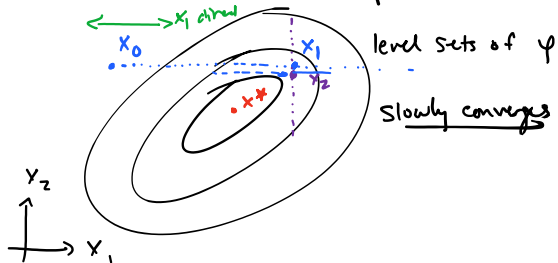
$\leftarrow \geq 0$ and $A \succ 0$ if $m > n$ + full rank

Set $\nabla \varphi(x) = 0$ and solve,

$$\nabla \varphi(x) = Ax - b, \text{ so find } x \text{ st } Ax = b$$

⚠ If you do want to solve least-squares, don't form $A = \tilde{A}^T \tilde{A}$. Instead use LSQR
CG-variant

One idea to min $\varphi(x)$ is coordinate descent / alternating minimization



Conjugate Directions

$\{p_i\}$ are conjugate directions if they are A-orthogonal

i.e., $\langle p_i, A p_j \rangle = 0$ if $i \neq j$ (Fact: $\langle x, y \rangle_A = x^T A y$ is an inner-product)

Notation: $\underbrace{\langle x | A | y \rangle}_{\langle x, y \rangle_A} = \langle x, A y \rangle = \langle A x, y \rangle = x^T A y = (A x)^T y \quad (A = A^T \text{ or } A = A^*)$

So, $\langle p_i | A | p_j \rangle = 0$ if $i \neq j$

If we have $\{p_i\}_{i=0}^{n-1}$ A -orthog, it's a basis.

Ex: p_i are eigenvectors: if A is symm, its eigenvectors are orthog.

$$\underbrace{\langle p_i | A | p_j \rangle}_{\lambda_j p_j} = \lambda_j \langle p_i | p_j \rangle = 0 \text{ if } i \neq j$$

Goal: find $\{p_i\}$ more cheaply than eigenvectors

Conjugate Direction method (abstract): Assume $\{p_i\}_{i=0}^{n-1}$ are conjugate directions.

$$(*) \quad x_{k+1} = x_k + \alpha_k p_k \quad (\text{like coordinate descent})$$

α_k solves $\min_{\alpha} \varphi(x_k + \alpha p_k)$ (exact linesearch)

\swarrow closed form $\varphi(x) = \frac{1}{2} \langle x | A | x \rangle - \langle b | x \rangle$

$$\alpha_k = - \frac{\langle A x_k - b | p_k \rangle}{\langle p_k | A | p_k \rangle} \quad r_k := A x_k - b$$

Thm: $x_n = x^*$, i.e., the true minimizer.

proof: $\{p_i\}$ is a basis, $x^* - x_0 = \sum_{i=0}^{n-1} \sigma_i p_i$ \swarrow coeff.

$$p_k^T A (x^* - x_0) = \sum_{i=0}^{n-1} \sigma_i \langle p_k | A | p_i \rangle = \sigma_k \langle p_k | A | p_k \rangle$$

$$\text{so } \sigma_k = \frac{\langle p_k | A | x^* - x_0 \rangle}{\langle p_k | A | p_k \rangle} \quad \forall k=0, 1, \dots, n-1$$

Also,

$$x_k = x_0 + \alpha_0 p_0 + \alpha_1 p_1 + \dots + \alpha_{k-1} p_{k-1} \quad \text{via } (*) \text{ recursively}$$

$$x_k - x_0 = \sum_{i=0}^{k-1} \alpha_i p_i$$

$$p_k^T A (x_k - x_0) = \sum_{i=0}^{k-1} \alpha_i \langle p_k | A | p_i \rangle = 0$$

$$\text{so } \langle p_k | A | x_k - x_0 \rangle = 0$$

$$A x^* = b$$

$$\text{claim } \underline{\underline{r_k = \frac{\langle p_k | A | x^* - x_k \rangle}{\langle p_k | A | p_k \rangle} = \frac{\langle p_k | b - Ax_k \rangle}{\dots} = \underline{\underline{\alpha_k}}}}$$

$$\Rightarrow x_n = x^* \quad \square$$

Facts via any CD method,

- $r_{k+1} = r_k + \alpha_k A p_k$
- $\langle r_k, p_i \rangle = 0 \quad i < k$
- x_k minimizes φ over $\text{K}(r_0, k-1) = \text{Span} \{ r_0, A r_0, \dots, A^{k-1} r_0 \}$
↖ Krylov subspace
 ie., $\text{K}(r_0, n-1) = \mathbb{R}^n$
- $\langle r_k, r_i \rangle = 0 \quad i < k$
- $p_k, r_k \in \text{K}(r_0, k)$

Conjugate Gradient = C.D. method that cheaply compute p_k

- x_0 anything, $r_0 = Ax_0 - b$, $p_0 = -r_0$

- Iterate

$$p_k = -r_k + \beta_k \cdot p_{k-1}, \quad \beta_k = \frac{\langle r_k | A | p_{k-1} \rangle}{\langle p_{k-1} | A | p_{k-1} \rangle}$$

$$x_{k+1} = x_k + \alpha_k p_k \quad \text{chosen so } \langle p_k | A | p_{k-1} \rangle = 0$$

$$\alpha_k = \frac{-\langle r_k | p_k \rangle}{\langle p_k | A | p_k \rangle}$$

Magic : $\langle p_k | A | p_i \rangle = 0 \quad \forall i \leq k-1$ ($i=k-1$ is by design)
 $i < k$ "are for free"
 (see Nocedal + Wright's text)

Cost: one matrix-vector multiply per step

Convergence Result (NW book)

$$\|x_k - x^*\|_A \leq 2 \cdot \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|x_0 - x^*\|_A$$

$\kappa = \kappa(A)$ is the condition # = $\frac{\max \text{ sing. value}}{\min \text{ sing. value}}$

Practical issues exist (eg., stability/roundoff), not discussed here.

• Nonlinear Conjugate Gradient

Linear CG: $\min_x \varphi(x) := \frac{1}{2} \langle x | A | x \rangle - \langle b | x \rangle$, $\nabla \varphi(x) = Ax - b = r$
 i.e., solve $Ax = b$ ← hence name "linear"

Nonlinear CG: $\min_x \varphi(x)$, $\varphi(x)$ isn't quadratic, $\nabla \varphi$ isn't linear

Linear:

$$x_{k+1} = x_k + \alpha_k p_k, \quad \alpha_k = \underset{\alpha}{\operatorname{argmin}} \varphi(x_k + \alpha p_k) \text{ known in closed form}$$

$$p_{k+1} = \underbrace{-r_{k+1}}_{=\nabla \varphi(x_{k+1})} + \beta_{k+1} p_k, \quad \beta_{k+1} = \frac{\langle r_{k+1} | A | p_k \rangle}{\langle p_k | A | p_k \rangle}$$

Nonlinear: $x_{k+1} = x_k + \alpha_k p_k$, $\alpha_k \approx \underset{\alpha}{\operatorname{argmin}} \varphi(x_k + \alpha p_k)$ Not in closed form.

$$p_{k+1} = \underbrace{-\nabla \varphi(x_{k+1})}_{\nabla \varphi_{k+1} \text{ for short}} + \beta_{k+1} p_k$$

Approximate via linesearch.
 ⚠ Sensitive! Need accurate linesearch

} Somewhat like Nesterov / momentum methods

β_{k+1} has many choices, all of which reduce to $\frac{\langle r_{k+1} | A | p_k \rangle}{\langle p_k | A | p_k \rangle}$ if φ quadratic

$$\text{Ex: } \beta_{k+1}^{FR} = \frac{\|\nabla \varphi_{k+1}\|^2}{\|\nabla \varphi_k\|^2}$$

others too (Polyak-Ribiere...)

Comments on non-linear CG

- can be fast, but limited to unconstrained optimization

Doesn't play well w/ constraints, even simple $x \geq 0$ ones

→ i.e., if φ is quadratic, CG is essentially optimal among all 1st order methods. So if "almost quadratic", non-linear CG might do very well

- Finicky w/ α, β terms. Hager + Zhang have many advanced versions, some require many (eg. 20) parameters.

- Global convergence theory isn't great

- Nemirovsky + Yudin proved there exist reasonable convex functions for which non-linear CG is slower than gradient descent

- IMO, quasi-Newton methods about as fast, and more stable (simpler)

(ie., if you tuned all non-linear CG param. just right, it might be faster, but parameters are problem dependent!)