

Lemma 2 Suppose \bar{x} minimizes f over $x^0 + \text{span}(s^1 \dots s^k)$ and s^{k+1} is A -conjugate to each s^i . Then

$$\hat{x} = \underset{x = \bar{x} + \alpha s^{k+1}}{\text{argmin}} f(x) \text{ is the minimizer of } f \text{ over } \text{span}(s^1, \dots, s^{k+1}) + x_0.$$

Proof. Essential same as previous separation we have seen. \square

This gives the Conjugate Gradient Method

Given $x_0 \in \mathbb{R}^d$, $s^0 = r^0 = b - Ax^0 = -\nabla f(x^0)$.

Iterate:

$$\alpha_i = \underset{\alpha}{\text{argmin}} f(x^i + \alpha s^i) \quad \leftarrow \alpha_i = \frac{s^{iT}(b - Ax^i)}{\langle s^i, s^i \rangle_A}$$

$$x^{i+1} = x^i + \alpha_i s^i$$

$$r^{i+1} = -\nabla f(x^{i+1}) = b - Ax^{i+1}$$

$$s^{i+1} = r^{i+1} - \sum_{j=1}^i \frac{\langle r^{i+1}, s^j \rangle_A}{\langle s^j, s^j \rangle_A} s^j \quad \leftarrow \text{for all } j \neq i$$

$$= r^{i+1} - \frac{\langle r^{i+1}, s^i \rangle_A}{\langle s^i, s^i \rangle_A} s^i$$

Theorem The Conjugate Gradient Method has

1. $\text{span}(r^0 \dots r^k) = \text{span}(s^0 \dots, s^k)$ has dimension $k+1$.
2. x^{k+1} minimizes f over $x^0 + \text{span}(r^0 \dots, r^k)$.

Proof. 1. is Gram-Schmidt using Lemma 1 for independence.

2. is what Lemma 2 tells us. \square

Claim. For $j < i$, $\langle r^{i+1}, s^j \rangle_A = 0$.

Proof. Let $L = \text{span}(r^0, \dots, r^i) = \text{span}(s^0, \dots, s^i)$

Theorem ensures x^{i+1} minimizes f over $x^0 + L$.

By Lemma 1, $-\nabla f(x^{i+1}) = r^{i+1}$ is orthogonal to L .

$$\Rightarrow r^{i+1}{}^T r^j = 0 \quad \forall j \leq i$$

$$\begin{aligned} \text{Then } \langle r^{i+1}, s^j \rangle_A &= r^{i+1}{}^T A s^j \\ &= r^{i+1}{}^T A (x^{j+1} - x^j) \\ &= r^{i+1}{}^T ((b - Ax^j) - (b - Ax^{j+1})) \\ &= r^{i+1}{}^T (r^j - r^{j+1}) \\ &= \underbrace{r^{i+1}{}^T r^j}_{=0} - \underbrace{r^{i+1}{}^T r^{j+1}}_{=0} \\ &\quad \text{if } j < i. \quad \square \end{aligned}$$

4. Convergence Rates

Recall Gradient Descent, we had $\delta_k = f(x_k) - f(x^*)$

$$\delta_k \leq \left(1 - \frac{1}{\text{Cond}(A)}\right)^k \delta_0$$

$$\text{where } \text{Cond}(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \left(= \frac{L}{\mu}\right)$$

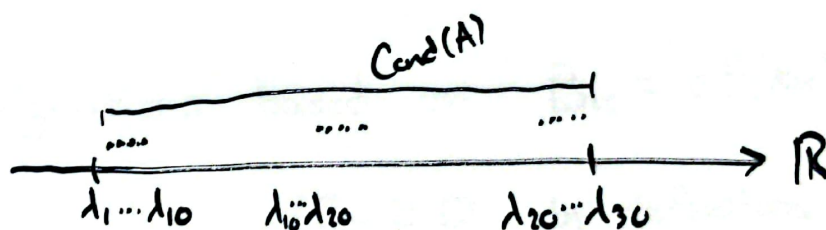
Conjugate Gradient does better (achieves optimal rate)

Theorem CGM has

$$\begin{aligned} \delta_k &\leq \left(\frac{\sqrt{\text{Cond}(A)} - 1}{\sqrt{\text{Cond}(A)} + 1} \right)^k \delta_0 \\ &\leq \left(1 - \frac{1}{\sqrt{\text{Cond}(A)}} \right)^k \delta_0. \end{aligned}$$

Proof. Spectral analysis of A . □

Faster guarantees for special arrangements of $\lambda_1, \dots, \lambda_n$



GMRES, Krylov Subspace find x_k minimizing over

$$K_k = \{b, Ab, A^2b, \dots, A^{k-1}b\}$$

Conjugate Gradient (More Generally)

(Caley-Hamilton Theorem
 $\Rightarrow A^{-1}b \in K_n$)

min $f(x)$, $x_{k+1} = x_k - \alpha_k s_k$

$$s_{k+1} = -\nabla f(x_{k+1}) + \beta_k s_k$$

$$\beta_k = \frac{\nabla f(x_k)^T (\nabla f(x_k) - \nabla f(x_{k-1}))}{s_k^T (\nabla f(x_k) - \nabla f(x_{k-1}))}$$

Nonlinear Least Squares

$$F(x) = 0 : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$\min f(x) = \frac{1}{2} \|F(x)\|_2^2$$

Recall HW 4, $\nabla f(x) = \nabla F(x)^T F(x)$

$$\nabla^2 f(x) = \underbrace{\nabla F(x)^T \nabla F(x)}_{\text{pretty nice}} + \underbrace{\sum_{i=1}^d \nabla^2 F_i(x) F_i(x)}_{\text{terrible}}$$

In general, f is nonconvex, best we can do

$$\nabla f(x) = \nabla F(x)^T F(x) = 0$$

Two Algorithms based on $B_k = \nabla F(x_k)^T \nabla F(x_k)$.

($B_k \succeq 0$ by definition)

Lets assume $B_k \succ 0$, ∇F be ind col).

($B_k \rightarrow \nabla^2 f(x^*)$ if $x_k \rightarrow x^*$)

1. Gauss-Newton Method

2. Levenberg-Marquardt Method (Trust-Region).

Each step of Gauss-Newton solves

$$\begin{aligned} P_k &= \operatorname{argmin} \quad \nabla f(x_k)^T p + \frac{1}{2} p^T B_k p \\ &= \operatorname{argmin} \quad \underline{F(x_k)^T \nabla F(x_k) p} + \frac{1}{2} p^T \underline{\nabla F(x_k)^T \nabla F(x_k) p}. \end{aligned}$$

$$\begin{aligned} \text{Solved by } \nabla F(x_k)^T \nabla F(x_k) p &= -\nabla F(x_k)^T F(x_k), \\ & (B_k p = -\nabla f(x_k)) \end{aligned}$$

$$\begin{aligned} &= \operatorname{argmin} \quad \frac{1}{2} \left\| \underline{F(x_k) + \nabla F(x_k) p} \right\|_2^2 \\ & \quad \text{Linearized } F(x_k + p) \end{aligned}$$

At each step we solve linear least squares problem.

(Conjugate Gradient Method good here)

If B_k has eigenvalues uniformly bounded away from 0 and ∞ , then we have good descent. Old results $O(1/\sqrt{k})$ convergence.

If $x_k \rightarrow x^*$, then $B_k \rightarrow \nabla^2 f(x^*)$, then $P_k \rightarrow$ "Newton Step" $\rightarrow \nabla^2 f(x^*)^{-1} \cdot \nabla f(x_k)$

(Superlinear rate)

Levenberg-Marquardt

$$S_k = \operatorname{argmin} \frac{1}{2} \| F(x_k) + \nabla F(x_k)s \|^2_2 \quad (1)$$

$\nwarrow m_k(s)$

$$\text{s.t. } \|s_k\|_2 \leq \delta_k$$

Well-defined always, always descent (for small δ_k)

How to pick δ_k ?

How to solve (1)?

(Trust Regions
after break)