

Lecture 23

Scribe?

Last time

- ▷ L-BFGS
- ▷ Conjugate gradient method

Today

- ▷ CG continued
- ▷ Convergence Guarantees
- ▷ Nonlinear least squares

Recall from last class

Lemma ⚡: Let x_0 and s_1, \dots, s_k be any vectors. Consider x_{k+1} given by $(*)$, then $\nabla f(x_{k+1})$ is orthogonal (in the standard sense) to $\text{span}\{s_1, \dots, s_k\}$.

Proof: Equivalently

$$y^* \in \arg \min_{y \in \mathbb{R}^k} f(x_0 + Sy)$$

By 1st-order optimality conditions:

$$S^T \nabla f(\overset{x_{k+1}}{x_0 + S y^*}) = 0$$

$\Rightarrow \nabla f(x_{k+1})$ is orthogonal to $\text{span}\{s_1, \dots, s_k\}$. \square

Thanks to separability:

Lemma \star : Suppose that x_{k+1} is given by (\star) and s_{k+1} is A -conjugate to each s_i . Then,

$$x_{k+2} \in \underset{\text{s.t. } x = x_{k+1} + \text{span}\{s_{k+1}\}}{\text{argmin}} f(x)$$

is also a solution of

$$x_{k+2} \in \underset{\text{s.t. } x = x_0 + \text{span}\{s_1, \dots, s_{k+1}\}}{\text{argmin}} f(x)$$

CG Method

Input: $x_0 \in \mathbb{R}^d$, $s_0 = r_0 = b - Ax_0$

Update $i \leq d$:

$$\alpha_i = \underset{\alpha}{\text{argmin}} f(x_i + \alpha s_i) \leftarrow$$

$$\alpha_i = \frac{s_i^T (b - Ax_i)}{\langle s_i, s_i \rangle_A}$$

$$x_{i+1} = x_i + \alpha_i s_i$$

$$r_{i+1} = -\nabla f(x_{i+1}) = b - Ax_{i+1}$$

$$s_{i+1} = r_{i+1} - \sum P_{s_i}^A(r_{i+1}) \leftarrow \begin{array}{l} \text{Lemma below} \\ r_{i+1} - P_{s_i}^A(r_{i+1}) \\ \text{Gram-Schmidt} \end{array}$$

Theorem: The conjugate gradient method has

1. $\text{span}\{r_1, \dots, r_k\} = \text{span}\{s_0, \dots, s_k\}$.

2. x_{k+1} is given by (★). +

Proof

1. Gram-Schmidt + Lemma \Rightarrow for independence.

2. Given by Lemma \square . \square

CG simplifies a lot:

Lemma: For $j < i$, $\langle r_{i+1}, s_j \rangle_A = 0$.

Proof: Let $L = \text{span}\{r_0, \dots, r^i\} = \text{span}\{s^0, \dots, s^i\}$.

The Theorem ensures that

x_{i+1} minimizes f over $x_0 + L$.

\Rightarrow By Lemma \Rightarrow , $-\nabla f(x_{i+1}) = r_{i+1}$ is orthogonal to L .

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

Expanding

$$\langle r_{i+1}, s_j \rangle_A = r_{i+1}^T A s_j.$$

$$= \frac{1}{\alpha_j} r_{i+1}^T A (\gamma_{j+1} - x_j)$$

$$= \frac{1}{\alpha_j} r_{i+1}^T ((b - Ax_j) - (b - Ax_{j+1}))$$

$$= \frac{1}{\alpha_i} \left(\underbrace{r_{i+1}^T r_j}_{=0} - \underbrace{r_{i+1}^T r_{j+1}}_{=0} \right) \leftarrow \text{By Lemma 3}$$

□

This ensures that we don't need to make a lot of unnecessary matrix-vector multiplies.

Convergence guarantees

Recall that with GD we had

$$f(x_k) - \min f \leq \left(1 - \frac{1}{\kappa(A)}\right)^k (f(x_0) - \min f)$$

where $\kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} = \frac{L}{\mu}$

You proved in a HW.

AGD achieved a faster convergence rate with $\sqrt{\kappa(A)}$ instead of $\text{cond}(A)$.

CG does just as well (it's optimal)

Theorem: The iterates of CG satisfy

$$f(x_k) - \min f \leq \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^k (f(x_0) - \min f)$$

$$\leq \left(1 - \frac{1}{\sqrt{\kappa(A)}}\right)^k (f(x_0) - \min f).$$

→

We are not going to prove this result, as the proof involves some matrix analysis and uses Chebychev polynomials.

Remarks:

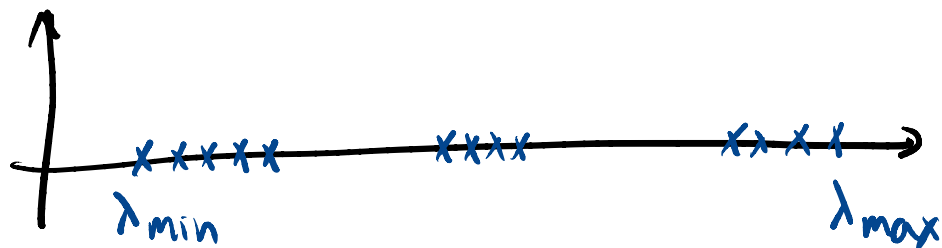
▷ The convergence is way better if $\kappa(A) \approx 1$. A natural idea is to precondition:

$$PAP^T \overset{\text{Invertible}}{y} = Pb$$

$\Rightarrow x = P^T y$ is a solution of $Ax = b$.

Active research area: How to come up with good preconditioners?

▷ For linear systems CG is often preferred over AGD. One reason is it offers faster convergence when eigenvalues are clustered, e.g.,



▷ How about asymmetric A ?

GMRES is a popular algorithm that updates

$$x_{k+1} = \arg \min \frac{1}{2} \|Ax - b\|^2$$

Krylov subspace. s.t. $x \in x_0 + L_k$ ← Computed via Arnoldi.

where $L_{k+1} = \text{span} \{r_0, Ar_0, A^2r_0, \dots, A^k r_0\}$.

By the Cayley-Hamilton Theorem $A^{-1}b \in L_n$.

This was invented by Saad and Schultz in 1986.

Krylov subspace methods (CG, GMRES, ...) are one of the Top 10 algorithms of the past century (according to SIAM).

▷ How about extensions for general f ?

They exist, but the guarantees and performance are not as strong;

see Chapter 5.2 of Nocedal & Wright.

Non linear least squares

Assume we have a mapping $r : \mathbb{R}^d \rightarrow \mathbb{R}^n$ and our goal is to minimize

$$\min_x f(x) = \frac{1}{2} \|r(x)\|_2^2$$

From HW4, we have

$$\nabla f(x) = \nabla r(x)^T r(x)$$

$$\nabla^2 f(x) = \underbrace{\nabla r(x)^T \nabla r(x)}_{\text{cheap to compute}} + \underbrace{\sum_{i=1}^n \nabla^2 r_i(x) r_i(x)}_{\text{expensive}}$$

small near a solution
↓

Our goal is to find a first-order critical point.

Two of the most popular first-order methods are

- ▷ Gauss-Newton
- ▷ Levenberg-Marguaret Method.

Gauss-Newton method

Similarly to Newton pick a direction
via

Well-defined if $B_k > 0$

$$p_k = \arg \min_p \nabla f(x_k)^T p + \frac{1}{2} p^T \underbrace{\nabla r(x_k)^T \nabla r(x_k)}_{B_k} p$$

\Leftrightarrow

$$\nabla r(x_k)^T \nabla r(x_k) p_k = -\nabla r(x_k)^T r(x_k)$$

\Leftrightarrow

$$p_k = \operatorname{argmin} \frac{1}{2} \underbrace{\| r(x_k) + \nabla r(x_k) p \|_2^2}_{\text{linearization of } r(x_k+p)}$$

We will not show it, but if

$$\mu I \leq B_k \leq L I$$

Then, this method ^{combined with appropriate step sizes} has descent and globally converges.

When x_k is close to $x^* \Rightarrow B_k$ is close to $\nabla^2 f(x_k)$ and the method has superlinear convergence.

(Chapter 10 of Nocedal & Wright).

Question: What can we do when B_k is not positive definite?

Levenberg - Marquardt Method

Idea: Bypass the lack of unique solutions by adding a norm constraint:

$$(\because) p_{k+1} = \operatorname{argmin} \frac{1}{2} \| r(x_k) + \nabla r(x_k)^T p \|^2$$

s.t. $\| p \| \leq \Delta_k.$

Trust-region

This prevents the need to pick α_k , but forces to pick Δ_k .

Q: How do we pick Δ_k ?

Q: How do we solve (\because) ?

We will cover trust-region after the break.