
Lecture 5: Second-order methods: Newton's method for unconstrained optimization (continued)

Coralia Cartis, Mathematical Institute, University of Oxford

C6.2/B2: Continuous Optimization

Global convergence of linesearch-Newton's method

- recall backtracking Armijo (bArmijo) linesearch.

Theorem 9 Let $f \in \mathcal{C}^2(\mathbb{R}^n)$ be bounded below on \mathbb{R}^n .

Let ∇f be Lipschitz continuous. Let the eigenvalues of $\nabla^2 f(x^k)$ be positive and uniformly bounded below, away from zero (for all k). Apply Newton's method to f with bArmijo linesearch and $\epsilon = 0$. Then

either

there exists $l \geq 0$ such that $\nabla f(x^l) = 0$

or

$\|\nabla f(x^k)\| \rightarrow 0$ as $k \rightarrow \infty$. \square

- Theorem 9 is satisfied if $f \in \mathcal{C}^2$ with ∇f Lipschitz continuous is also strongly convex (i.e., the eigenvalues of $\nabla^2 f(x)$ for all x are positive, bounded below, away from zero). Then s^k is

descent for all k .

[Much stronger conditions than for SD methods.]

Global convergence of linesearch-Newton's method ...

Proof of Theorem 9. The conditions of Theorem 4 (Global convergence of GLM with bArmijo linesearch) are satisfied. Thus Th 4 gives that either $\exists l \geq 0$ such that $\nabla f(x^l) = 0$ or

$$M_k := \min \left\{ \frac{|\nabla f(x^k)^T s^k|}{\|s^k\|}, |\nabla f(x^k)^T s^k| \right\} \longrightarrow 0 \text{ as } k \rightarrow \infty. \quad (\dagger)$$

Let $\nabla^2 f(x^k) := H_k$. Th assumptions on $f \implies \forall s \in \mathbb{R}^n, s \neq 0,$

$$0 < \lambda_{\min} \leq \lambda_{\min}(H_k) \leq \frac{s^T H_k s}{\|s\|^2} \leq \lambda_{\max}(H_k) \leq \lambda_{\max}.$$

$$\begin{aligned} |\nabla f(x^k)^T s^k| &= |\nabla f(x^k)^T H_k^{-1} \nabla f(x^k)| \geq \lambda_{\min}(H_k^{-1}) \|\nabla f(x^k)\|^2 \\ &= \frac{\|\nabla f(x^k)\|^2}{\lambda_{\max}(H_k)} \geq \frac{\|\nabla f(x^k)\|^2}{\lambda_{\max}}. \end{aligned}$$

$$\|s^k\|^2 = \nabla f(x^k)^T H_k^{-2} \nabla f(x^k) \leq \lambda_{\max}(H_k^{-2}) \|\nabla f(x^k)\|^2 \leq \lambda_{\min}^{-2} \|\nabla f(x^k)\|^2.$$

$$\implies M_k \geq \min \left\{ \frac{\lambda_{\min}}{\lambda_{\max}} \|\nabla f(x^k)\|, \frac{1}{\lambda_{\max}} \|\nabla f(x^k)\|^2 \right\} \text{ for all } k$$

$$\implies \nabla f(x^k) \longrightarrow 0 \text{ as } k \rightarrow \infty. \quad \square$$

Global convergence for general second-order GLMs

In GLM, let s^k be defined by $B^k s^k = -\nabla f(x^k)$, where B^k symmetric and positive definite matrix.

Theorem 10 Let $f \in \mathcal{C}^1(\mathbb{R}^n)$ be bounded below on \mathbb{R}^n .

Let ∇f be Lipschitz continuous. Let the eigenvalues of B^k be uniformly bounded above and below, away from zero, for all k .

Apply GLM with above s^k and bArmijo linesearch and $\epsilon = 0$.

Then

either

there exists $l \geq 0$ such that $\nabla f(x^l) = 0$

or

$\|\nabla f(x^k)\| \rightarrow 0$ as $k \rightarrow \infty$. \square

- Theorem requires locally **strongly convex** quadratic models of f for all k (but the Hessian of f may not be pos. def.).

Modified damped Newton methods

If $\nabla^2 f(x^k)$ is not positive definite, it is usual to solve instead

$$\left(\nabla^2 f(x^k) + M^k\right) s^k = -\nabla f(x^k),$$

where

- M^k chosen such that $\nabla^2 f(x^k) + M^k$ is “sufficiently” positive definite.
- $M^k := 0$ when $\nabla^2 f(x^k)$ is “sufficiently” positive definite.

Options:

1. As $\nabla^2 f(x^k)$ is symmetric, we can factor $\nabla^2 f(x^k) = Q^k D^k (Q^k)^\top$, where Q^k is orthogonal and D^k is diagonal, and set

$$\nabla^2 f(x^k) + M^k := Q^k \max(\epsilon I, |D^k|) (Q^k)^\top,$$

for some “small” $\epsilon > 0$. Expensive approach for large problems.

Modified damped Newton methods

2. Estimate $\lambda_{\min}(\nabla^2 f(x^k))$ and set

$$M^k := \max(0, \epsilon - \lambda_{\min}(\nabla^2 f(x^k)))I.$$

Cheaper. Often tried in practice but “biased” (may overemphasize a large negative eigval at the expense of small, positive ones).

3. Modified Cholesky: compute Cholesky factorization

$$\nabla^2 f(x^k) = L^k (L^k)^\top,$$

where L^k is lower triangular matrix. Modify the generated L^k if the factorization is in danger of failing (modify small or negative diagonal pivots, etc.).

Popular in computations.

Approximating the Hessian matrix by finite differences

Approximating the Hessian from gradient vals: $i \in \{1, \dots, n\}$;

$$[\nabla^2 f(x)]e^i \approx \frac{1}{h}[\nabla f(x + he^i) - \nabla f(x)]$$

Cost of approximating $\nabla^2 f(x)$ is $n + 1$ gradient values.

For all finite-differencing, careful with the choice of h in computations:

- “too large” $h \rightarrow$ inaccurate approximations,
- “too small” $h \rightarrow$ numerical cancellation errors.

But successful techniques exist for smooth noiseless problems when sufficient function and/or gradient values can be computed.

For noisy problems, use **derivative-free optimization** methods (if problem size is not too large).

Quasi-Newton methods

Secant approximations for computing $B^k \approx \nabla^2 f(x^k)$

At the start of the GLM, choose B^0 (say, $B^0 := I$). After computing $s^k = -(B^k)^{-1} \nabla f(x^k)$ and $x^{k+1} = x^k + \alpha^k s^k$, compute update B^{k+1} of B^k .

Wish list:

Compute B^{k+1} as a function of already-computed quantities $\nabla f(x^{k+1}), \nabla f(x^k), \dots, \nabla f(x^0), B^k, s^k$,

B^{k+1} should be symmetric, nonsingular (pos. def.),

B^{k+1} “close” to B^k , a “cheap” update of B^k , $B^k \rightarrow \nabla^2 f(x^k)$, etc.

\implies a new class of methods: faster than steepest descent method, cheaper to compute per iteration than Newton's.

For the first wish, choose B^{k+1} to satisfy the secant equation

$$\gamma^k := \nabla f(x^{k+1}) - \nabla f(x^k) = B^{k+1}(x^{k+1} - x^k) = B^{k+1} \alpha^k s^k.$$

Quasi-Newton methods ...

Interpretation of the secant equation:

It is satisfied by $B^{k+1} := \nabla^2 f$ when f is a quadratic function.

The change in gradient contains information about the Hessian.

The gradient change predicted by the current quadratic model

$$\nabla f(x^{k+1}) - \nabla f(x^k) \approx \nabla q(x^k + \alpha^k s^k) - \nabla q(x^k) = -\alpha^k \nabla f(x^k),$$

where $q(x^k + s) = f(x^k) + \nabla f(x^k)^\top s + \frac{1}{2} s^\top B^k s$

and $s^k = -(B^k)^{-1} \nabla f(x^k)$.

Want the new quadratic model

$$u(x^k + s) := f(x^k) + \nabla f(x^k)^\top s + \frac{1}{2} s^\top B^{k+1} s$$

to predict correctly the change in gradient γ^k , i.e.,

$$\gamma^k = \nabla f(x^{k+1}) - \nabla f(x^k) = \nabla u(x^{k+1}) - \nabla u(x^k) = B^{k+1} (x^{k+1} - x^k).$$

Quasi-Newton methods ...

Many ways to compute B^{k+1} to satisfy the secant equation.
Trade-off between “wishes” on the list for some of the methods.

Symmetric rank 1 updates.

[see Prob Sheet 3]

Set $B^{k+1} := B^k + u^k (u^k)^\top$, for some $u^k \in \mathbb{R}^n$, and all $k \geq 0$.

- B^{k+1} symmetric, “close” to B^k .
- Work per iteration: $\mathcal{O}(n^2)$ (as opposed to the $\mathcal{O}(n^3)$ of Newton), due to Sherman-Morrison-Woodbury formula!

The secant equation $\implies u^k = (\gamma^k - B^k \delta^k) / \rho^k$,
where $\delta^k := x^{k+1} - x^k = \alpha^k s^k$, $(\rho^k)^2 := (\gamma^k - B^k \delta^k)^\top \delta^k > 0$.

- B^k may not be positive definite, s^k may not be descent.
- ρ^k may be close to zero leading to large updates.

Other updates: **BFGS**, **DFP**, Broyden family, etc.

Quasi-Newton methods ...

BFGS updates.

[see Prob Sheet 3]

- Broyden-Fletcher-Goldfarb-Shanno (independently).

Set $B_{k+1} := B_k + u_k u_k^\top + v_k v_k^\top$, for some $u_k \in \mathbb{R}^n$, $v_k \in \mathbb{R}^n$.

- It is a rank 2 update (if u_k and v_k are linearly independent).
- SWM formula yields $\mathcal{O}(n^2)$ operations/iteration.
- In practice, update the Cholesky factors of B_k (still $\mathcal{O}(n^2)$).

Given $B_k = J_k J_k^\top$, where J_k arbitrary nonsingular, and $\|\cdot\|_F$ Frobenius norm, let J_{k+1} solve

$$\min_J \|J - J_k\|_F \quad \text{subject to} \quad J \delta_k = \gamma_k.$$

$$\Rightarrow B_{k+1} := J_{k+1} J_{k+1}^\top = B_k + u_k u_k^\top + v_k v_k^\top,$$

where $u_k u_k^\top = -B_k \delta_k \delta_k^\top B_k / (\delta_k^\top B_k \delta_k)$, $v_k v_k^\top = \gamma_k \gamma_k^\top / (\gamma_k^\top \delta_k)$.

- Let $J_k := L_k$ the lower triangular Cholesky factor of B_k .
-

Quasi-Newton methods ...

BFGS updates. (continued)

- Thus B_{k+1} is “close” to B_k .
 - B_k symmetric pos. def. \Rightarrow B_{k+1} symmetric **pos. def.** (provided $(\delta^k)^T \gamma^k > 0$, ensured by say, **Wolfe linesearch**)
 - **BFGS method:** GLM with $s_k := -B_k^{-1} \nabla f(x_k)$, with B_k updated by BFGS formula on each iteration.
 - For global convergence of BFGS method, must use **Wolfe linesearch** to compute stepsize instead of bArmijo linesearch.
 - The BFGS method has **local Q-superlinear convergence!**
 - When applying the BFGS method with exact linesearches, to a strictly convex quadratic function f , then $B_k = \nabla^2 f$ after n iterations.
 - Satisfies all the wishes on the wish list! Has been very popular when second derivatives of f are not available.
-

Appendix: providing derivatives to algorithms

How to compute/provide derivatives to a solver?

- Calculate derivatives by hand when easy/simple objective and constraints; user provides code that computes them.
- Calculate or approximate derivatives automatically:
 - Automatic differentiation: breaks down computer code for evaluating f into elementary arithmetic operations + differentiate by chain rule. Software: ADIFOR, ADOL-C.
 - Symbolic differentiation: manipulate the algebraic expression of f (if available). Software: symbolic packages of MAPLE, MATHEMATICA, MATLAB.
 - Finite differencing \longrightarrow approximate derivatives.

See Nocedal & Wright, Numerical Optimization (2nd edition, 2006) for more details of the above procedures.