

Math 273a: Optimization
The Barzilai-Borwein method

Instructor: Wotao Yin
Department of Mathematics, UCLA
Fall 2015

Main features of the Barzilai-Borwein (BB) method

- The BB method was published in a 8-page paper¹ in 1988
- It is a gradient method with special step sizes. The method is motivated by Newton's method but does not compute Hessian
- At nearly no extra cost over the standard gradient method, the method is often found to significantly outperform the standard gradient method
- The method is used along with non-monotone line search as a convergence safeguard for non-quadratic problems

¹J. Barzilai and J. Borwein. Two-point step size gradient method. IMA J. Numerical Analysis 8, 141–148, 1988.

Background

Goal: minimize $f(\mathbf{x})$, where f is a smooth function

Let $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ and $\mathbf{F}^{(k)} = \nabla^2 f(\mathbf{x}^{(k)})$.

- **gradient method:** $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$
 - choice of α_k : fixed, exact line search, or backtracking line search
 - **pros:** simple
 - **cons:** no use of 2nd order information, relatively slow progress
- **Newton's method:** $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{F}^{(k)})^{-1} \mathbf{g}^{(k)}$
 - **pros:** 2nd-order information, 1-step for quadratic function, fast convergence near solution
 - **cons:** forming and computing $(\mathbf{F}^{(k)})^{-1}$ is expensive, need modifications if $\mathbf{F}^{(k)} \neq 0$
- **BB method:** choose α_k so that $\alpha_k \mathbf{g}^{(k)}$ “approximates” $(\mathbf{F}^{(k)})^{-1} \mathbf{g}^{(k)}$

Derive the BB method

- Consider quadratic optimization

$$\underset{\mathbf{x}}{\text{minimize}} \quad q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x},$$

where $\mathbf{A} \succ 0$ is symmetric. Gradient is $\mathbf{g}^{(k)} = \mathbf{A} \mathbf{x}^{(k)} - \mathbf{b}$. Hessian is \mathbf{A} .

- Newton step: $\mathbf{d}_{\text{newton}}^{(k)} = -\mathbf{A}^{-1} \mathbf{g}^{(k)}$
- Goal:** choose α_k so that $-\alpha_k \mathbf{g}^{(k)} = -(\alpha_k^{-1} \mathbf{I})^{-1} \mathbf{g}^{(k)}$ approximates $\mathbf{d}_{\text{newton}}^{(k)}$
- Define: $\mathbf{s}^{(k-1)} := \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}$ and $\mathbf{y}^{(k-1)} := \mathbf{g}^{(k)} - \mathbf{g}^{(k-1)}$. \mathbf{A} satisfies:

$$\mathbf{A} \mathbf{s}^{(k-1)} = \mathbf{y}^{(k-1)}.$$

- Therefore, given $\mathbf{s}^{(k-1)}$ and $\mathbf{y}^{(k-1)}$, how about choose α_k so that

$$(\alpha_k^{-1} \mathbf{I}) \mathbf{s}^{(k-1)} \approx \mathbf{y}^{(k-1)}$$

- **Goal:**

$$(\alpha_k^{-1} I) \mathbf{s}^{(k-1)} \approx \mathbf{y}^{(k-1)}.$$

- **BB method:**

- Least-squares problem: (let $\beta = \alpha^{-1}$)

$$\alpha_k^{-1} = \arg \min_{\beta} \frac{1}{2} \|\mathbf{s}^{(k-1)} \beta - \mathbf{y}^{(k-1)}\|^2 \implies \alpha_k^1 = \frac{(\mathbf{s}^{(k-1)})^T \mathbf{s}^{(k-1)}}{(\mathbf{s}^{(k-1)})^T \mathbf{y}^{(k-1)}}$$

- Alternative Least-squares problem:

$$\alpha_k = \arg \min_{\alpha} \frac{1}{2} \|\mathbf{s}^{(k-1)} - \mathbf{y}^{(k-1)} \alpha\|^2 \implies \alpha_k^2 = \frac{(\mathbf{s}^{(k-1)})^T \mathbf{y}^{(k-1)}}{(\mathbf{y}^{(k-1)})^T \mathbf{y}^{(k-1)}}$$

- α_k^1 and α_k^2 are called the BB step sizes.

Apply the BB method

- At $k = 0$, $\mathbf{x}^{(k-1)}$ and $\mathbf{g}^{(k-1)}$ (and thus $\mathbf{s}^{(k-1)}$ and $\mathbf{y}^{(k-1)}$) are unavailable, so apply 1 iteration of the standard gradient descent.
- Then, switch to the BB method at $k = 1$
- We can use either α_k^1 or α_k^2 for all $k \geq 1$, or alternate between them
- We can also fix $\alpha_k = \alpha_k^1$ or $\alpha_k = \alpha_k^2$ for a few consecutive steps and then alternate.
- It performs very well on minimizing both quadratic and other differentiable functions
- However, f_k and $\|\nabla f_k\|$ are **not** monotonic!

Numerical: steepest descent vs BB on quadratic programming

- **Model:**

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}.$$

- The template of a gradient iteration

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^{(k)} - \alpha_k (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{b}).$$

- **Steepest descent** selects $\alpha_k = \arg \min_{\alpha} f(\mathbf{x}^{(k)} - \alpha_k (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{b}))$, so

$$\alpha_k = \frac{(\mathbf{r}^k)^T \mathbf{r}^{(k)}}{(\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^{(k)}}$$

where $\mathbf{r}^{(k)} := \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}$.

- **BB** selects α_k as

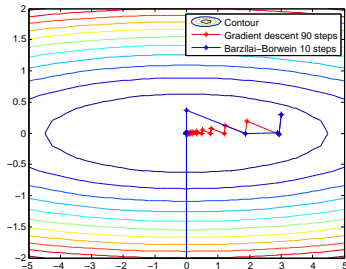
$$\alpha_k = \frac{(\mathbf{s}^{(k-1)})^T \mathbf{s}^{(k-1)}}{(\mathbf{s}^{(k-1)})^T \mathbf{y}^{(k-1)}}$$

Numerical example

- Set symmetric matrix \mathbf{A} to have the condition number $\frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})} = 50$.
- Stopping criterion:

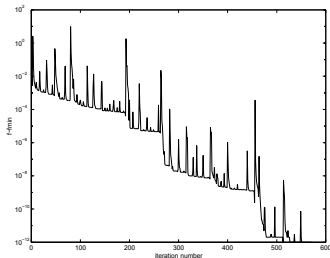
$$\|\mathbf{r}^{(k)}\| < 10^{-8}$$

- **Steepest descent** took 90 iterations to stop
- **BB** took only 10 iterations to stop (went very far temporarily and then came back)



Properties of Barzilai-Borwein

- For quadratic functions, it has R-linear convergence²
- For 2D quadratic function, it has Q-superlinear convergence³
- No convergence guarantee for smooth convex problems. On these problems, we pair up BB with non-monotone line search.



$$\text{BB on Laplace2: } \min \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x} + \frac{h^2}{4} \sum_{ijk} u_{ijk}^4.$$

²Dai and Liao [2002]

³Barzilai and Borwein [1988], Dai [2013]

Safeguard: nonmonotone line search

- Definition: line search that permits temporary growth but enforces overall descent of the function value
- For nonconvex problems, they improve the likelihood of global optimality
- Improve convergence speed when a monotone scheme is forced to creep along the bottom of a narrow curved valley
- Early nonmonotone line search method⁴ developed for Newton's methods

$$f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) \leq \max_{0 \leq j \leq m_k} f(\mathbf{x}^{k-j}) + c_1 \alpha \nabla f_k^T \mathbf{d}^{(k)}$$

However, it may still kill R-linear convergence. **Example:** $x \in \mathbb{R}$,

$$\underset{x}{\text{minimize}} f(x) = \frac{1}{2}x^2, \quad x^0 \neq 0, \quad \mathbf{d}^{(k)} = -x^{(k)}.$$

$$\alpha_k = \begin{cases} 1 - 2^{-k}, & k = i^2 \text{ for some integer } i, \\ 2, & \text{otherwise,} \end{cases}$$

converges R-linear but fails to satisfy the condition for k large.

⁴Grippo, Lampariello, and Lucidi [1986]

Zhang-Hager nonmonotone line search⁵

1. initialize $0 < c_1 < c_2 < 1$, $C_0 \leftarrow f(\mathbf{x}^0)$, $Q_0 \leftarrow 1$, $\eta < 1$, $k \leftarrow 0$
2. while *not converged* do
 - 3a. compute α_k satisfying the modified Wolfe conditions OR
 - 3b. find α_k by backtracking, to satisfy the modified Armijo condition:
sufficient decrease: $f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) \leq C_k + c_1 \alpha_k \nabla f_k^T \mathbf{d}^{(k)}$
4. $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$
5. $Q_{k+1} \leftarrow \eta Q_k + 1$, $C_{k+1} \leftarrow (\eta Q_k C_k + f(\mathbf{x}^{k+1})) / Q_{k+1}$.

Comments:

- If $\eta = 1$, then $C_k = \frac{1}{k+1} \sum_{j=0}^k f_j$.
- Since $\eta < 1$, C_k is a weighted sum of all past f_j , more weights on recent f_j .

⁵Zhang and Hager [2004]

Convergence (advanced topic)

The results below are left to the reader as an exercise.

If $f \in C^1$ and bounded below, $\nabla f_k^T \mathbf{d}^{(k)} < 0$, then

- $f_k \leq C_k \leq \frac{1}{k+1} \sum_{j=0}^{(k)} f_j$
- there exists α_k satisfying the modified Wolfe or Armijo conditions

In addition, if ∇f is Lipschitz with constant L , then

- $\alpha_k > C \frac{|\nabla f_k^T \mathbf{d}^{(k)}|}{\|\mathbf{d}^{(k)}\|}$ for some constant depending on c_1, c_2, L and the backing factor

Furthermore, if for all sufficiently large k , we have uniform bounds

$$\nabla f_k^T \mathbf{d}^{(k)} \leq -c_3 \|\nabla f_k\|^2 \quad \text{and} \quad \|\mathbf{d}^{(k)}\| \leq c_4 \|\nabla f_k\|$$

then ▪ $\lim_{k \rightarrow \infty} \nabla f_k = 0$

Once again, pairing with non-monotone linear search, Barzilai-Borwein gradient methods *work every well on general unconstrained differentiable problems.*

References:

- Yu-Hong Dai and Li-Zhi Liao. R-linear convergence of the Barzilai and Borwein gradient method. *IMA Journal of Numerical Analysis*, 22(1):1–10, 2002.
- J. Barzilai and J.M. Borwein. Two-point step size gradient methods. *IMA Journal of Numerical Analysis*, 8(1):141–148, 1988.
- Yu-Hong Dai. A new analysis on the barzilai-borwein gradient method. *Journal of the Operations Research Society of China*, pages 1–12, 2013.
- Luigi Grippo, Francesco Lampariello, and Stephano Lucidi. A nonmonotone line search technique for Newton's method. *SIAM Journal on Numerical Analysis*, 23(4): 707–716, 1986.
- Hongchao Zhang and William W Hager. A nonmonotone line search technique and its application to unconstrained optimization. *SIAM Journal on Optimization*, 14(4): 1043–1056, 2004.