

# The Barzilai-Borwein method

Chenglong Bao

YMSC, Tsinghua University

Acknowledgement: slides based on Prof. Wotao Yin(UCLA).

## Main features of the Barzilai-Borwein (BB) method

$$X_{k+1} = X_k - \underbrace{(\nabla^2 f(x_k))^{-1}}_{\text{Newton.}} \nabla f(x_k)$$

- The BB method was published in a 8-page paper<sup>1</sup> 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

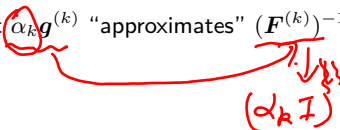
---

<sup>1</sup>J. 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  $\alpha_k$ : 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  $\alpha_k$  so that  $\alpha_k \mathbf{g}^{(k)}$  "approximates"  $(\mathbf{F}^{(k)})^{-1} \mathbf{g}^{(k)}$   


# Derive the BB method

$$f(x)=0$$

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

- Consider quadratic optimization

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

$$f'(x_k) = \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$$

$$\Rightarrow f'(x_k)(x_k - x_{k-1}) = f(x_k) - f(x_{k-1})$$

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

- Newton step:  $d_{\text{newton}}^{(k)} = -A^{-1}g^{(k)}$

$\approx A'$

- Goal:** choose  $\alpha_k$  so that  $-\alpha_k g^{(k)} = -(\alpha_k^{-1} I)^{-1} g^{(k)}$  approximates  $d_{\text{newton}}^{(k)}$

- Define:  $\underline{s^{(k-1)} := x^{(k)} - x^{(k-1)}}$  and  $\underline{y^{(k-1)} := g^{(k)} - g^{(k-1)}}$ .  $A$  satisfies:

$$\underline{As^{(k-1)} = y^{(k-1)}}. \quad \leftarrow \underline{\text{secant equation.}}$$

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

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

- Goal:

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

- BB method:

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

$$\alpha_k^{-1} = \arg \min_{\beta} \frac{1}{2} \| \underline{s^{(k-1)} \beta - y^{(k-1)}} \|^2 \Rightarrow \alpha_k^1 = \frac{(s^{(k-1)})^T s^{(k-1)}}{\underline{(s^{(k-1)})^T y^{(k-1)}}} \approx 0.$$

- Alternative Least-squares problem: (BB-2)

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

- $\alpha_k^1$  and  $\alpha_k^2$  are called the BB step sizes.

$$\langle y^{(k-1)}, s^{(k-1)} \rangle \geq 0 \Leftrightarrow \begin{aligned} y^{(k-1)} &= \nabla f(x^{(k)}) - \nabla f(x^{(k-1)}) \\ s^{(k-1)} &= x^{(k)} - x^{(k-1)} \end{aligned}$$

## 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  $\alpha_k^1$  or  $\alpha_k^2$  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:**

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

- 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  $\alpha_k$  as

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

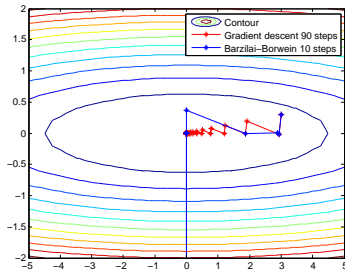
## Numerical example

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

- Stopping criterion:

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

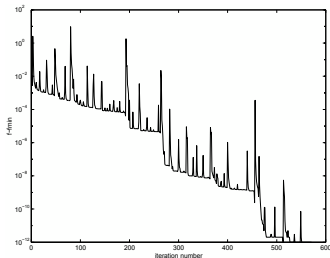
- Steepest descent** took 90 iterations to stop
- BB** took only 10 iterations to stop (went very far temporarily and then came back)





# Properties of Barzilai-Borwein

- For quadratic functions, it has R-linear convergence<sup>2</sup>
- For 2D quadratic function, it has Q-superlinear convergence<sup>3</sup>
- No convergence guarantee for smooth convex problems. On these problems, we pair up BB with non-monotone line search.



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

---

<sup>2</sup>Dai and Liao [2002]

<sup>3</sup>Barzilai 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 method<sup>4</sup> developed for Newton's methods



$$d^{(k)T} \nabla f(x^k) < 0$$

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

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

$$x_{k+1} = \begin{cases} -x_k, & \text{if } k \neq i^2 \\ \frac{1}{2^{-k}} x_k, & \text{if } k = i^2 \end{cases} \quad \text{minimize } f(x) = \frac{1}{2} x^2, \quad x^0 \neq 0, \quad d^{(k)} = -x^{(k)}.$$

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

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

<sup>4</sup>Grippo, Lampariello, and Lucidi [1986]

## Zhang-Hager nonmonotone line search<sup>5</sup>

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  $\alpha_k$  satisfying the modified Wolfe conditions OR
  - 3b. find  $\alpha_k$  by backtracking, to satisfy the modified Armijo condition:

sufficient decrease:  $f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) \leq \underline{\underline{C_k}} + c_1 \alpha_k \nabla f_k^T \mathbf{d}^{(k)}$  ✓

4.  $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$

5.  $\boxed{Q_{k+1} \leftarrow \eta Q_k + 1, C_{k+1} \leftarrow (\eta Q_k C_k + f(\mathbf{x}^{k+1}))/Q_{k+1}.}$

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  $f_j$ .

<sup>5</sup>Zhang 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  $\alpha_k$  satisfying the modified Wolfe or Armijo conditions

In addition, if  $\nabla f$  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

$$\mathbf{d}^k = -\nabla f(x^k)$$

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

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

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.