

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{s.t.} & x \in E \end{array}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a real-valued function and E is the feasible region.

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$. If $f \in C^1$ then

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \frac{\partial f}{\partial x_2}(x) \\ \vdots \\ \frac{\partial f}{\partial x_r}(x) \\ \vdots \\ \frac{\partial f}{\partial x_{n-1}}(x) \\ \frac{\partial f}{\partial x_n}(x) \end{bmatrix}$$

- Derivative of f at a point $x := Df(x) = [\nabla f(x)]^T$.
- sometimes we will use $\nabla f(x) = g(x)$ and $\nabla f(x_k) = g(x_k) = g_k$

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$. If $f \in C^2$ then Hessian of $f(x)$ at point x is $H(x) = \nabla g(x)^T = \nabla \{\nabla^T f(x)\}$. Hence

$$H(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(x) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(x) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x) & \frac{\partial^2 f}{\partial x_2^2}(x) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(x) & \frac{\partial^2 f}{\partial x_n \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(x) \end{bmatrix}$$

- $H(x)$ is an $n \times n$ square symmetric matrix.

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. If $f \in C^1$ then Jacobian of $f(x)$ at point x is $J(x) = Df(x)$ and is defined as

$$J(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(x) & \frac{\partial f_1}{\partial x_2}(x) & \cdots & \frac{\partial f_1}{\partial x_n}(x) \\ \frac{\partial f_2}{\partial x_1}(x) & \frac{\partial f_2}{\partial x_2}(x) & \cdots & \frac{\partial f_2}{\partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(x) & \frac{\partial f_m}{\partial x_2}(x) & \cdots & \frac{\partial f_m}{\partial x_n}(x) \end{bmatrix}$$

- Here $x = [x_1, x_2, x_3, \dots, x_n]^T$ and $f(x) = [f_1(x), f_2(x), f_3(x), \dots, f_m(x)]^T$
- $J(x)$ is an $m \times n$ matrix.

- Let $g : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$.

Let $f : (a, b) \rightarrow \Omega$.

Define $h : (a, b) \rightarrow \mathbb{R}$ by $h(t) = g(f(t))$

Then $Dh(t) = \frac{dh}{dt} = [Dg(f(t))][Df(t)]$, i.e. $= \langle Dg(f(t)), [Df(t)] \rangle$
 $= \left[\frac{\partial g}{\partial x_1}(f(t)) \quad \frac{\partial g}{\partial x_2}(f(t)) \quad \cdots \quad \frac{\partial g}{\partial x_n}(f(t)) \right] \cdot \left[\frac{\partial f_1}{\partial t} \quad \frac{\partial f_2}{\partial t} \quad \cdots \quad \frac{\partial f_n}{\partial t} \right]^T$

- Let $f, g : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

Define $h : \mathbb{R}^n \rightarrow \mathbb{R}$ by $h(x) = [f(x)]^T g(x)$

Then $Dh(x) = [f(x)]^T Dg(x) + [g(x)]^T Df(x)$. Particularly

- $D(y^T Ax) = y^T A$
- $D(x^T Ax) = x^T A + x^T A^T = 2x^T A$ (if A is symmetric)
- $D(y^T x) = y^T$
- $D(x^T x) = x^T + x^T = 2x^T$

Definition

The level set of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at level c is the set $\{x \in \mathbb{R}^n : f(x) = c\}$.

Particularly

- for $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, it is a curve.
- for $f : \mathbb{R}^3 \rightarrow \mathbb{R}$, it is a surface.

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

- Directional derivative of f at a point x in the direction of d is $\langle g(x), d \rangle$.

Problem:

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{s.t.} & x \in \mathbb{R}^n \end{array}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a real-valued function.

- Directional derivative of f at a point x in the direction of d is $\langle g(x), d \rangle$.
- The gradient g and its negative $-g$ are the steepest ascent and steepest descent directions.

Input: Function $f(x)$ and initial guess x^0 .

Until Convergence Do:

iteration:

$$x^{k+1} = x^k + \alpha_k d(x^k) \quad (1)$$

where x^k is the current estimate of a local minimizer to the problem.
 $d(x^k)$: is the current search direction in the \mathbb{R}^n . In above equation:

$$d(x^k) = -M_k g^k \quad (2)$$

where $g^k = \nabla f(x^k)$ and M_k is $n \times n$ matrix. And, α_k (step size) is the minimizer of a function $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$ defined as

$$\begin{aligned} \phi(\alpha) &= f(x^k + \alpha d(x^k)) \\ &= f(x^k - \alpha M_k g^k) \end{aligned}$$

- $g^k = 0$ (In computation gradient is rarely identically zero).
- $\|g^k\| < \varepsilon$ for some pre specified small positive value ε .
- $\|f(x^{k+1}) - f(x^k)\| < \varepsilon$ or $\|x^{k+1} - x^k\| < \varepsilon$.
- $\frac{\|f(x^{k+1}) - f(x^k)\|}{\|f(x^k)\|} < \varepsilon$ or $\frac{\|x^{k+1} - x^k\|}{\|x^k\|} < \varepsilon$ (relative stopping criteria that is scale-independent).
- $\frac{\|f(x^{k+1}) - f(x^k)\|}{\max\{1, \|f(x^k)\|\}} < \varepsilon$ or $\frac{\|x^{k+1} - x^k\|}{\max\{1, \|x^k\|\}} < \varepsilon$ (to avoid division by a small number).

Output: $x^* = x^{k+1}$ and $f^* = f(x^*)$

$$x^{k+1} = x^k + \alpha_k d(x^k)$$

where

$$d(x^k) = -M_k g^k$$

and, α_k (step size) is the minimizer of a function $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$ defined as

$$\begin{aligned}\phi(\alpha) &= f(x^k + \alpha d(x^k)) \\ &= f(x^k - \alpha M_k g^k)\end{aligned}$$

Different choice of matrix M_k gives a different name to our algorithm.

- $M_k = I_n$: Steepest Descent Method (SD).
- If M_k is generated by conjugate vectors with respect to an appropriate PD matrix: Conjugate Gradient Method (CG), Fletcher Reeves Method (FR).
- $M_k = [H(x^k)]^{-1} = H_k^{-1}$: Newton Method.
- If M_k is some approximation of Hessian inverse: Quasi-Newton Type Method.

$$x^{k+1} = x^k + \alpha_k d(x^k)$$

where

$$d(x^k) = -M_k g^k$$

and, α_k (step size) is the minimizer of a function $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$ defined as

$$\begin{aligned}\phi(\alpha) &= f(x^k + \alpha d(x^k)) \\ &= f(x^k - \alpha M_k g^k)\end{aligned}$$

Definition

Any iterative method as defined in the main frame-work is said to have descent property if $f(x^{k+1}) < f(x^k)$, $\forall k$, provided $g^k \neq 0$.

Definition

Any iterative method as defined in the main frame-work is said to have quadratic termination property if the minimum of $f(x) = \frac{1}{2}x^T A x - b^T x + C$, A is $n \times n$ SPD matrix, is reached in at most n iteration, $b \in \mathbb{R}^n$, C is a constant.

Definition

Any iterative method as defined in the main frame-work is said to be globally convergent if $x^k \rightarrow x^$ (a local minimizer of f) for any initial guess x^0 as $k \rightarrow \infty$.*

Definition

Any iterative method as defined in the main frame-work is said to have order of convergence p if there exists $0 < a < \infty$ such that

$$\lim_{k \rightarrow \infty} \frac{\|x^{k+1} - x^*\|}{\|x^k - x^*\|^p} = a$$

If $p = 1$ and $a = 0$ then it is called super linear convergent method.

minimize $f(x)$ s.t. $x \in \mathbb{R}^n$

✓ Input $f(x)$, x^0 , and ε (for stopping condition)

✓ **Until Convergence Do:** for $k = 0, 1, 2, 3, \dots$

Step 1 Calculate $g^k := \nabla f(x^k)$

Step 2 Set $d^k = -g^k$

Step 3 Find α_k , the value of α that minimizes $f(x^k + \alpha d^k)$

Step 4 Set $x^{k+1} = x^k + \alpha_k d^k$

✓ EndDo

Theorem

SD is a descent method.

Proof:

$$x^{k+1} = x^k + \alpha_k d^k$$

where $d(x^k) = -g^k$ and, α_k is the minimizer of a function $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$ defined as $\phi(\alpha) = f(x^k + \alpha d^k)$. Since

$$\begin{aligned} \frac{d\phi}{d\alpha} \Big|_{\alpha=0} &= [Df(x^k + \alpha d(x^k)) \Big|_{\alpha=0}] d^k \\ &= [g^k]^T d^k = -[d^k]^T d^k = -\|d^k\|^2 < 0 \end{aligned}$$

So there exists $\bar{\alpha} > 0$ such that $\phi(\alpha) < \phi(0) \ \forall \ 0 < \alpha \leq \bar{\alpha}$.

Hence, we have

$$\phi(\alpha_k) = f(x^k + \alpha_k d^k) = f(x^{k+1}) \leq \phi(\alpha) < \phi(0) = f(x^k).$$

Problem: minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

Step 1 Calculate $g^k := \nabla f(x^k)$.

Here: $g^k = Ax^k - b = r^k$ (say)

Step 2 Set $d^k = -M_k g^k$. **Here:** $d^k = -M_k r^k$.

Step 3 Find α_k , the value of α that minimizes $\phi(\alpha) := f(x^k + \alpha M_k d^k)$.

Here: $\alpha_k = -\frac{\langle M_k d^k, r^k \rangle}{\langle M_k d^k, A M_k d^k \rangle}$

Step 4 Set $x^{k+1} = x^k + \alpha_k M_k d^k$.

Here: $x^{k+1} = x^k - \frac{\langle M_k d^k, r^k \rangle}{\langle M_k d^k, A M_k d^k \rangle} M_k d^k$

Problem: minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

Step 1 Calculate $g^k := \nabla f(x^k)$. **Here:** $g^k = Ax^k - b = r^k$ (say)

Step 2 Set $d^k = -g^k$. **Here:** $d^k = -r^k$.

Step 3 Find α_k , the value of α that minimizes $\phi(\alpha) := f(x^k + \alpha d^k)$.

Here: $\alpha_k = \frac{\langle r^k, r^k \rangle}{\langle r^k, Ar^k \rangle}$

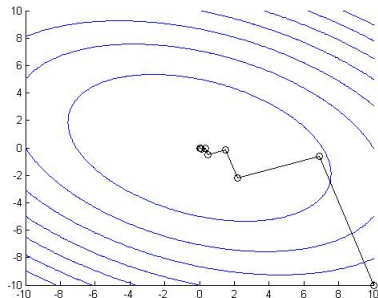
Step 4 Set $x^{k+1} = x^k + \alpha_k d^k$. **Here:** $x^{k+1} = x^k - \frac{\langle r^k, r^k \rangle}{\langle r^k, Ar^k \rangle} r^k$

SD does not have quadratic termination property

Problem: minimize $x_1^2 + x_1x_2 + 2x_2^2$.

Here: minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, $b = 0$ and $A = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}$.

<i>Iter</i>	x_1	x_2
1	10.0000	-10.0000
2	6.8749	-0.6247
3	2.1875	-2.1875
10	0.0157	-0.0014
12	0.0034	-0.0003
15	0.0002	-0.0002
18	0.0000	-0.0000



Theorem

SD moves in orthogonal directions

Proof:

$$x^{k+1} = x^k + \alpha_k d^k, \quad x^{k+2} = x^{k+1} + \alpha_{k+1} d^{k+1}$$

So,

$$\langle x^{k+2} - x^{k+1}, x^{k+1} - x^k \rangle = \alpha_k \alpha_{k+1} \langle d^{k+1}, d^k \rangle \quad (3)$$

Since, α_k is the minimizer of a function $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$ defined as $\phi(\alpha) = f(x^k + \alpha d^k)$. So

$$\begin{aligned} \left. \frac{d\phi}{d\alpha} \right|_{\alpha=\alpha_k} &= [Df(x^k + \alpha d(x^k))]|_{\alpha=\alpha_k} d^k \\ &= [g^{k+1}]^T d^k = -[d^{k+1}]^T d^k = 0 \end{aligned} \quad (4)$$

Above two equations imply the conclusion.

Theorem

SD is a globally convergent method.

SD convergence rate is linear for quadratic functions

For simplicity, we take

Problem: minimize $f(x) = \frac{1}{2}x^T A x$, A is SPD. It is known here that $x^* = 0$, $f^* = f(x^*) = 0$. Let max and min eigenvalues of matrix A are B and b respectively.

We know that here SD generates the following sequence for minimizer:

$x^{k+1} = x^k - \alpha_k r^k$ where $\alpha_k = \frac{\langle r^k, r^k \rangle}{\langle r^k, A r^k \rangle}$ and $r^k = \nabla f(x^k) = A x^k$.

We can show easily that $\frac{f(x^{k+1})}{f(x^k)} = 1 - \beta$ where $\beta = \frac{\langle r^k, r^k \rangle^2}{\langle r^k, A r^k \rangle \langle r^k, A^{-1} r^k \rangle}$.

By Rayleigh inequality:

$\frac{b}{2} \|x^{k+1}\|^2 \leq f(x^{k+1})$ and $\frac{B}{2} \|x^k\|^2 \geq f(x^k)$. Hence

$\frac{b}{2} \|x^{k+1}\|^2 \leq f(x^{k+1}) = (1 - \beta) f(x^k) \leq (1 - \beta) \frac{B}{2} \|x^k\|^2$

$\Rightarrow \frac{\|x^{k+1}\|^2}{\|x^k\|^2} = \frac{B}{b} (1 - \beta) \leq \frac{B}{b} \frac{B-b}{B}$ (By Kantorovich inequality). Hence

$$\frac{\|x^{k+1}\|}{\|x^k\|} \leq \sqrt{\frac{B-b}{b}}$$

$$\text{minimize } f(x) \text{ s.t. } x \in \mathbb{R}^n$$

✓ Input $f(x)$, x^0 , and ε (for stoping condition)

✓ **Until Convergence Do:** for $k = 0, 1, 2, 3, \dots$

Step 1 Calculate $g^k := \nabla f(x^k)$

Step 2 Set $d^k = -[H(x^k)]^{-1}g^k$

Step 3 Set $x^{k+1} = x^k + d^k$

✓ EndDo

Idea:

Let x^k be the current estimator of x^* . Take quadratic approximation q of f near x^k and find its minima instead of f . Minima of q is treated as the next estimator x^{k+1} for the x^* . Thus,

$q(x) = f(x^k) + g(x^k)^T(x - x^k) + \frac{1}{2}(x - x^k)^T H(x^k)(x - x^k)$. This gives:

$$Dq(x) = g(x^k)^T + \frac{1}{2}((x - x^k)^T H(x^k) + (x - x^k)^T H(x^k)^T)$$

$$= g(x^k)^T + (x - x^k)^T H(x^k). \text{ This implies:}$$

$$\nabla q(x) = g(x^k) + H(x^k)(x - x^k) = 0 \text{ gives the step 3.}$$

minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

Until Convergence Do:

Step 1 Calculate $g^k := \nabla f(x^k)$

Step 2 Set $d^k = -[H(x^k)]^{-1}g^k$

Step 3 Set $x^{k+1} = x^k + d^k$

EndDo

First iteration:

Step 1 $g^0 := \nabla f(x^0) = A(x^0) - b$

Step 2 $d^0 = -A^{-1}g^0$
 $= -A^{-1}(A(x^0) - b)$
 $= -x^0 + A^{-1}b$

Step 3 $x^1 = x^0 - x^0 + A^{-1}b$
 $= A^{-1}b$

✓ $g^1 = A(x^1) - b = 0$ STOP

Newton's Method give exact answer in only one iteration always!

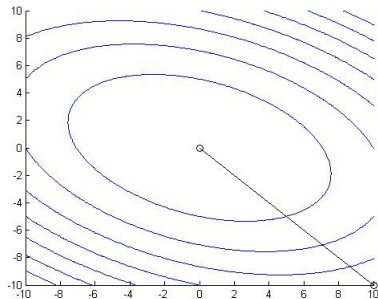
Why?

quadratic Approximation of a quadratic function.

Problem: minimize $x_1^2 + x_1x_2 + 2x_2^2$.

Here: minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, $b = 0$ and $A = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}$.

$$\begin{pmatrix} \text{Iter} & x_1 & x_2 \\ 0 & 10.0000 & -10.0000 \\ 1 & 0 & 0 \end{pmatrix}$$



(Any nonlinear function: not quadratic)

$$\text{minimize } f(x) \quad \text{s.t. } x \in \mathbb{R}^n$$

✓ Input $f(x)$, x^0 , and ε (for stoping condition)

✓ **Until Convergence Do:** for $k = 0, 1, 2, 3, \dots$

Step 1 Calculate $g^k := \nabla f(x^k)$

Step 2 Set $d^k = -[H(x^k)]^{-1}g^k$

Step 3 Set $x^{k+1} = x^k + d^k$

✓ EndDo

Problems:

- 1 $H(x^k)$ must be invertible for each iteration.
- 2 Above scheme is not descent for any arbitrary x_0 (Why?).

$$\text{minimize } f(x) \quad \text{s.t. } x \in \mathbb{R}^n$$

Modified Newton's Method:

- ✓ Input $f(x)$, x^0 , and ε (for stoping condition)
- ✓ **Until Convergence Do:** for $k = 0, 1, 2, 3, \dots$
 - Step 1** Calculate $g^k := \nabla f(x^k)$
 - Step 2** Set $d^k = -[H(x^k)]^{-1}g^k$ OR Solve $H(x^k)d^k = -g^k$.
 - Step 3** Find α_k , the value of α that minimizes $f(x^k + \alpha d^k)$
 - Step 4** Set $x^{k+1} = x^k + \alpha_k d^k$
- ✓ EndDo

Theorem

If Hessian is SPD in each iteration, then Modified Newton Method is descent.

Theorem

If Hessian is SPD in each iteration, then Modified Newton Method is descent.

Proof:

$$x^{k+1} = x^k + \alpha_k d^k$$

where $d(x^k) = -[H(x^k)]^{-1}g^k$ and, α_k is the minimizer of a function $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$ defined as $\phi(\alpha) = f(x^k + \alpha d^k)$. Since

$$\begin{aligned} \frac{d\phi}{d\alpha}|_{\alpha=0} &= [Df(x^k + \alpha d(x^k))|_{\alpha=0}]d^k \\ &= [g^k]^T d^k = [d^k]^T g^k = -[d^k]^T H(x^k) d^k < 0 \end{aligned}$$

So there exists $\bar{\alpha} > 0$ such that $\phi(\alpha) < \phi(0) \ \forall \ 0 < \alpha \leq \bar{\alpha}$.

Hence, we have

$$\phi(\alpha_k) = f(x^k + \alpha_k d^k) = f(x^{k+1}) \leq \phi(\alpha) < \phi(0) = f(x^k).$$

$$\text{minimize } f(x) \quad \text{s.t. } x \in \mathbb{R}^n$$

Levenberg-Marquardt Method:

- ✓ Input $f(x)$, x^0 , and ε (for stoping condition)
- ✓ **Until Convergence Do:** for $k = 0, 1, 2, 3, \dots$
 - Step 1** Calculate $g^k := \nabla f(x^k)$
 - Step 2** Solve $[H(x^k) + \mu_k I]d^k = -g^k$ for sufficiently large μ_k .
 - Step 3** Find α_k , the value of α that minimizes $f(x^k + \alpha d^k)$
 - Step 4** Set $x^{k+1} = x^k + \alpha_k d^k$
 - ✓ EndDo

$$\text{minimize } f(x) \quad \text{s.t. } x \in \mathbb{R}^n$$

Levenberg-Marquardt Method:

- ✓ Input $f(x)$, x^0 , and ε (for stopping condition)
- ✓ Calculate g^0 , $H(x^0)$, and set $\mu = 10^6$ (large number),
dif = TRUE.
- ✓ **while** $\|g^k\| < \varepsilon$
 - Step 1** if (dif) $\mu = 2\mu$; else $\mu = \frac{\mu}{2}$.
 - Step 2** Solve $[H(x^k) + \mu_k I]d^k = -g^k$.
 - Step 3** Set $x^{k+1} = x^k + d^k$
 - Step 4** Calculate f^{k+1} , g^{k+1} , H^{k+1}
 - Step 5** if $(f(x^{k+1}) < f(x^k))$ dif = FALSE; else dif = TRUE.

Definition

Let A be any SPD $n \times n$ matrix. Then Vectors $p^0, p^1, p^2, \dots, p^{n-1}$ are called A -conjugate directions iff $\langle p^i, Ap^j \rangle = 0, \forall i \neq j$.

- This definition generalize the orthogonality.
- Above n directions are LI in \mathbb{R}^n .

$$\text{minimize } f(x) \quad \text{s.t. } x \in \mathbb{R}^n$$

✓ Input $f(x)$, x^0 , n A -conjugate directions $p^0, p^1, p^2, \dots, p^{n-1}$ and ε (for stopping condition)

✓ **Until Convergence Do:** for $k = 0, 1, 2, 3, \dots$

Step 1 Find α_k , the value of α that minimizes $f(x^k + \alpha p^k)$

Step 2 Set $x^{k+1} = x^k + \alpha_k p^k$

✓ EndDo

Idea:

Instead of negative gradient direction, move in the conjugate directions.

minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

- ✓ Input $f(x)$, x^0 , n A -conjugate directions $p^0, p^1, p^2, \dots, p^{n-1}$ and ε (for stopping condition)
- ✓ **Until Convergence Do:** for $k = 0, 1, 2, 3, \dots$

Step 1 Find α_k , the value of α that minimizes $f(x^k + \alpha p^k)$.

Here: $\alpha_k = -\frac{\langle p^k, r^k \rangle}{\langle p^k, Ap^k \rangle}$, where $r^k = \nabla f(x^k) = Ax^k - b$.

Step 2 Set $x^{k+1} = x^k + \alpha_k p^k$. **Here:** $x^{k+1} = x^k - \frac{\langle p^k, r^k \rangle}{\langle p^k, Ap^k \rangle} p^k$.

✓ EndDo

Proof:

$\phi(\alpha) = f(x^k + \alpha p^k) = \frac{1}{2}(x^k + \alpha p^k)^T A(x^k + \alpha p^k) - b^T(x^k + \alpha p^k)$. Thus

$$\frac{d\phi}{d\alpha} = \frac{1}{2}[(p^k)^T A(x^k + \alpha p^k) + (x^k + \alpha p^k)^T A p^k] - b^T p^k$$

$$= (p^k)^T A x^k + \alpha (p^k)^T A p^k - (p^k)^T b.$$

$$\frac{d\phi}{d\alpha} = 0 \text{ gives } \alpha = \alpha_k.$$

Basic Conjugate Direction Method satisfies quadratic termination property

minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

✓ Input $f(x)$, x^0 , n A -conjugate directions $p^0, p^1, p^2, \dots, p^{n-1}$ and ε (for stopping condition)

Step 1 $\alpha_k = -\frac{\langle p^k, r^k \rangle}{\langle p^k, Ap^k \rangle}$, where $r^k = \nabla f(x^k) = Ax^k - b$.

Step 2 Set $x^{k+1} = x^k + \alpha_k p^k$.

Proof: By iterative step:

$$\begin{aligned}x^1 &= x^0 + \alpha_0 p^0 \\x^2 &= x^1 + \alpha_1 p^1 = x^0 + \alpha_0 p^0 + \alpha_1 p^1 \\x^k &= x^0 + \alpha_0 p^0 + \alpha_1 p^1 + \alpha_2 p^2 + \dots + \alpha_{k-1} p^{k-1} \quad (5)\end{aligned}$$

$$x^n = x^0 + \alpha_0 p^0 + \alpha_1 p^1 + \alpha_2 p^2 + \dots + \alpha_{n-1} p^{n-1} \quad (6)$$

We have to show that $x^n = x^*$. Since $x^* \in \mathbb{R}^n$ and $p^0, p^1, p^2, \dots, p^{n-1}$ is a basis of $x^* \in \mathbb{R}^n$, there exists λ_i 's such that

$$x^* - x^0 = \lambda_0 p^0 + \lambda_1 p^1 + \lambda_2 p^2 + \dots + \lambda_{n-1} p^{n-1} \quad (7)$$

From (6) and (7), it is clear that to prove $x^n = x^*$, it is sufficient to show that $\lambda_k = \alpha_k$.

Basic Conjugate Direction Method satisfies quadratic termination property

By premultiplying (7) with $(p^k)^T A$ we have

$$\lambda_k = \frac{\langle p^k, A(x^* - x^0) \rangle}{\langle p^k, Ap^k \rangle} \quad (8)$$

Now, $x^* - x^0 = (x^* - x^k) + (x^k - x^0)$ implies

$$\begin{aligned} \langle p^k, A(x^* - x^0) \rangle &= \langle p^k, A(x^* - x^k) \rangle + \langle p^k, A(x^k - x^0) \rangle \\ &= \langle p^k, A(x^* - x^k) \rangle \\ &\quad (\because \text{second term is zero by (5)}) \\ &= \langle p^k, b - Ax^k \rangle = -\langle p^k, r^k \rangle \end{aligned} \quad (9)$$

From (8) and (9), it is clear that $\lambda_k = \alpha_k$

One more important property of the Basic Conjugate Direction Method for quadratic problem

minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

✓ Input $f(x)$, x^0 , n A -conjugate directions $p^0, p^1, p^2, \dots, p^{n-1}$ and ε (for stopping condition)

Step 1 $\alpha_k = -\frac{\langle p^k, r^k \rangle}{\langle p^k, Ap^k \rangle}$, where $r^k = \nabla f(x^k) = Ax^k - b$.

Step 2 Set $x^{k+1} = x^k + \alpha_k p^k$.

Theorem : New residual \perp all old conjugate directions

For any $k \in \{0, 1, 2, \dots, n-1\}$, $\langle r^{k+1}, p^i \rangle = 0 \forall i = 0, 1, 2, \dots, k$.

Proof: For $k = 0$, see

$$\begin{aligned}\langle r^1, p^0 \rangle &= \langle Ax^1 - b, p^0 \rangle = \langle Ax^0 + \alpha_0 Ap^0 - b, p^0 \rangle \\ &= \langle Ax^0 - b + \alpha_0 Ap^0, p^0 \rangle = \langle r^0 + \alpha_0 Ap^0, p^0 \rangle \\ &= \langle r^0, p^0 \rangle + \alpha_0 \langle Ap^0, p^0 \rangle = 0.\end{aligned}$$

Last equality is due to the definition of α_0 .

One more important property of the Basic Conjugate Direction Method for quadratic problem

Now Assume that for some $k - 1$ and $i = 0, 1, 2, \dots, k - 1$ statement is correct, i.e., $\langle r^k, p^i \rangle = 0 \forall i = 0, 1, 2, \dots, k - 1$. Then for k and $i = 0, 1, 2, \dots, k - 1$, we have

$$\begin{aligned}\langle r^{k+1}, p^i \rangle &= \langle r^k + \alpha_k A p^k, p^i \rangle \\ &\quad \text{(by the iteration used in the algo: } r^{k+1} = r^k + \alpha_k A p^k \text{)} \\ &= \langle r^k, p^i \rangle + \alpha_k \langle A p^k, p^i \rangle = 0 \\ &\quad \text{(first term is zero by the assumption of induction)} \\ &\quad \text{(and second is zero by the conjugacy of } p^i \text{ vectors.)}\end{aligned}$$

Now its remains to show that $\langle r^{k+1}, p^k \rangle = 0$. To see it

$$\begin{aligned}\langle r^{k+1}, p^k \rangle &= \langle r^k + \alpha_k A p^k, p^k \rangle \\ &= \langle r^k, p^k \rangle + \alpha_k \langle A p^k, p^k \rangle = 0.\end{aligned}$$

Last equality is due to the definition of α_k .

Ques: How we find the A -Conjugate Directions

Problem in consideration: minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

Iteration step in Basic Conjugate Direction Method:

$x^{k+1} = x^k - \alpha_k p^k$, where $\alpha_k = \frac{\langle p^k, r^k \rangle}{\langle p^k, Ap^k \rangle}$ and $r^k = Ax^k - b$.

Construction of p^i , $i = 0, 1, 2, \dots, n-1$:

$$p^0 = -r^0 \quad (10)$$

$$p^k = -r^k + \beta_{k-1} p^{k-1}, \quad k = 1, 2, \dots, n-1 \quad (11)$$

where

$$\text{where } \beta_{k-1} = \frac{\langle r^k, Ap^{k-1} \rangle}{\langle p^{k-1}, Ap^{k-1} \rangle}, \quad k = 1, 2, \dots, n-1 \quad (12)$$

This choice of β_{k-1} in equation (11) gives the following facts:

Gradient vectors r^k , $k = 0, 1, 2, \dots, n-1$ are mutually orthogonal

For any $k \in \{0, 1, 2, \dots, n-1\}$, $\langle r^{k+1}, r^i \rangle = 0 \forall i = 0, 1, 2, \dots, k$.

p^k , $k = 0, 1, 2, \dots, n-1$ are mutually A conjugate directions

For any $k \in \{0, 1, 2, \dots, n-1\}$, $\langle p^{k+1}, Ap^i \rangle = 0 \forall i = 0, 1, 2, \dots, k$.

Theorem: Gradient vectors r^k , $k = 0, 1, 2, \dots, n-1$ are mutually orthogonal

Proof: for arbitrary i , using equation (11) we have

$$\begin{aligned} \langle r^{k+1}, r^i \rangle &= \langle r^{k+1}, -p^i + \beta_{i-1} p^{i-1} \rangle \\ &= -\langle r^{k+1}, p^i \rangle + \beta_{i-1} \langle r^{k+1}, p^{i-1} \rangle \\ &= 0, \end{aligned} \tag{13}$$

by the previous result that, for any $k \in \{0, 1, 2, \dots, n-1\}$, $\langle r^{k+1}, p^i \rangle = 0 \forall i = 0, 1, 2, \dots, k$.

Theorem: $p^k, k = 0, 1, 2, \dots, n-1$ are mutually A conjugate

Proof: p^1 and p^0 are A conjugate by the definition of β_0 . Now, we show that p^2 is A conjugate to p^1 and p^0 . It is clear that p^2 and p^1 are A conjugate by the definition of β_1 .

$$\begin{aligned} \langle p^2, Ap^0 \rangle &= \langle -r^2 + \beta_1 p^1, Ap^0 \rangle = -\langle r^2, Ap^0 \rangle + \beta_1 \langle p^1, Ap^0 \rangle \\ &= -\langle r^2, Ap^0 \rangle \quad (\because \text{second term } \langle p^1, Ap^0 \rangle = 0) \\ &= -\langle r^2, \frac{r^1 - r^0}{\alpha_0} \rangle = \frac{1}{\alpha_0} [\langle r^2, r^0 \rangle - \langle r^2, r^1 \rangle] \\ &\quad \because \text{iterative step } x^{k+1} = x^k + \alpha_k p^k \Rightarrow Ap^k = \frac{r^{k+1} - r^k}{\alpha_k} \\ &= 0 \text{ by the previous result.} \end{aligned}$$

Theorem: $p^k, k = 0, 1, 2, \dots, n-1$ are mutually A conjugate

Now Assume that for some $k-1$ and $i = 0, 1, 2, \dots, k-1$ statement is correct, i.e., $\langle p^k, p^i \rangle = 0 \forall i = 0, 1, 2, \dots, k-1$. Then for k and $i = 0, 1, 2, \dots, k-1$, we have

$$\begin{aligned}\langle p^{k+1}, Ap^i \rangle &= \langle -r^{k+1} + \beta_k p^k, Ap^i \rangle \\&= -\langle r^{k+1}, Ap^i \rangle + \beta_k \langle p^k, Ap^i \rangle \\&\quad (\text{second term is zero by the assumption of induction}) \\&= -\langle r^{k+1}, \frac{r^{i+1} - r^i}{\alpha_i} \rangle \\&= \frac{1}{\alpha_i} [\langle r^{k+1}, r^i \rangle - \langle r^{k+1}, r^{i+1} \rangle] \\&\quad \because \text{iterative step } x^{k+1} = x^k + \alpha_k p^k \Rightarrow Ap^k = \frac{r^{k+1} - r^k}{\alpha_k} \\&= 0 \text{ by the previous result.}\end{aligned}$$

minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, A is SPD.

Step 1 Set $k = 0$, and select the initial point x^0 .

Step 2 $g^0 = \nabla f(x^0) = Ax^0 - b = r^0$. If $g^0 = 0$ STOP else $p^0 = -g^0$.

Step 3 $\alpha_k = -\frac{\langle p^k, g^k \rangle}{\langle p^k, Ap^k \rangle}$.

Step 4 $x^{k+1} = x^k + \alpha_k p^k$

Step 5 $g^{k+1} = \nabla f(x^{k+1}) = Ax^{k+1} - b = r^{k+1}$. If $g^{k+1} = 0$ STOP.

Step 6 $\beta_k = \frac{\langle g^{k+1}, Ap^k \rangle}{\langle p^k, Ap^k \rangle}$

Step 7 $p^{k+1} = -g^{k+1} + \beta_k p^k$

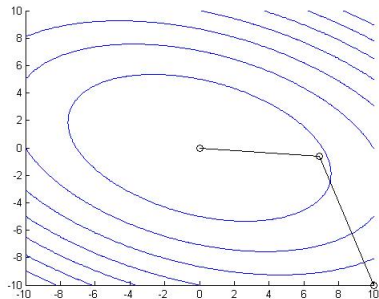
Step 8 Set $k = k + 1$ and go to Step 3.

Since CG method is a conjugate direction method so minimizes the given problem in at most n steps, where n is the order of the matrix A .

Problem: minimize $x_1^2 + x_1x_2 + 2x_2^2$.

Here: minimize $f(x) = \frac{1}{2}x^T Ax - b^T x$, $b = 0$ and $A = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}$.

<i>Iter</i>	x_1	x_2
0	10.0000	-10.0000
1	6.8750	-0.6250
2	0	0



minimize $f(x) : x \in \mathbb{R}^n$.

Step 1 Set $k = 0$, and select the initial point x^0 .

Step 2 $g^0 = \nabla f(x^0) = Ax^0 - b = r^0$. If $g^0 = 0$ STOP else $p^0 = -g^0$.

Step 3 $\alpha_k = -\frac{\langle p^k, g^k \rangle}{\langle p^k, Ap^k \rangle}$.

Step 4 $x^{k+1} = x^k + \alpha_k p^k$

Step 5 $g^{k+1} = \nabla f(x^{k+1}) = Ax^{k+1} - b = r^{k+1}$. If $g^{k+1} = 0$ STOP.

Step 6 $\beta_k = \frac{\langle g^{k+1}, Ap^k \rangle}{\langle p^k, Ap^k \rangle}$

Step 7 $p^{k+1} = -g^{k+1} + \beta_k p^k$

Step 8 Set $k = k + 1$ and go to Step 3.

minimize $f(x) : x \in \mathbb{R}^n$.

Step 1 Set $k = 0$, and select the initial point x^0 .

Step 2 $g^0 = \nabla f(x^0) = Ax^0 - b = r^0$. If $g^0 = 0$ STOP else
 $p^0 = -g^0$.

Step 3 $\alpha_k = -\frac{\langle p^k, g^k \rangle}{\langle p^k, Ap^k \rangle}$.

Step 4 $x^{k+1} = x^k + \alpha_k p^k$

Step 5 $g^{k+1} = \nabla f(x^{k+1}) = Ax^{k+1} - b = r^{k+1}$. If $g^{k+1} = 0$
 STOP.

Step 6 $\beta_k = \frac{\langle g^{k+1}, Ap^k \rangle}{\langle p^k, Ap^k \rangle}$.

Step 7 $p^{k+1} = -g^{k+1} + \beta_k p^k$

Step 8 Set $k = k + 1$ and go to Step 3.

minimize $f(x) : x \in \mathbb{R}^n$.

Step 1 Set $k = 0$, and select the initial point x^0 .

Step 2 $g^0 = \nabla f(x^0)$. If $g^0 = 0$ STOP else $p^0 = -g^0$.

Step 3 $\alpha_k = -\frac{\langle p^k, g^k \rangle}{\langle p^k, Ap^k \rangle}$.

Step 4 $x^{k+1} = x^k + \alpha_k p^k$

Step 5 $g^{k+1} = \nabla f(x^{k+1}) = Ax^{k+1} - b = r^{k+1}$. If $g^{k+1} = 0$
STOP.

Step 6 $\beta_k = \frac{\langle g^{k+1}, Ap^k \rangle}{\langle p^k, Ap^k \rangle}$.

Step 7 $p^{k+1} = -g^{k+1} + \beta_k p^k$

Step 8 Set $k = k + 1$ and go to Step 3.

minimize $f(x) : x \in \mathbb{R}^n$.

Step 1 Set $k = 0$, and select the initial point x^0 .

Step 2 $g^0 = \nabla f(x^0)$. If $g^0 = 0$ STOP else $p^0 = -g^0$.

Step 3 Find α_k , the value of α that minimizes $f(x^k + \alpha p^k)$.

Step 4 $x^{k+1} = x^k + \alpha_k p^k$

Step 5 $g^{k+1} = \nabla f(x^{k+1}) = Ax^{k+1} - b = r^{k+1}$. If $g^{k+1} = 0$
STOP.

Step 6 $\beta_k = \frac{\langle g^{k+1}, Ap^k \rangle}{\langle p^k, Ap^k \rangle}$.

Step 7 $p^{k+1} = -g^{k+1} + \beta_k p^k$

Step 8 Set $k = k + 1$ and go to Step 3.

minimize $f(x) : x \in \mathbb{R}^n$.

Step 1 Set $k = 0$, and select the initial point x^0 .

Step 2 $g^0 = \nabla f(x^0)$. If $g^0 = 0$ STOP else $p^0 = -g^0$.

Step 3 Find α_k , the value of α that minimizes $f(x^k + \alpha p^k)$.

Step 4 $x^{k+1} = x^k + \alpha_k p^k$

Step 5 $g^{k+1} = \nabla f(x^{k+1})$. If $g^{k+1} = 0$ STOP.

Step 6 $\beta_k = \frac{\langle g^{k+1}, A p^k \rangle}{\langle p^k, A p^k \rangle}$.

Step 7 $p^{k+1} = -g^{k+1} + \beta_k p^k$

Step 8 Set $k = k + 1$ and go to Step 3.

Recall: $\beta_k = \frac{\langle g^{k+1}, Ap^k \rangle}{\langle p^k, Ap^k \rangle}$

Actually replacement of A in the above formula is $H(x^k)$. But fortunately, algebraic manipulation in the above formula is possible with the knowledge of the function value $f(x^k)$ and gradient value g^k .

Hestenes-Stiefel (SF) Modification:

It is clear that iterative step that $x^{k+1} = x^k + \alpha_k p^k \Rightarrow Ap^k = \frac{g^{k+1} - g^k}{\alpha_k}$

Use this value in the formula of β_k in place of Ap^k , thus:

$$\beta_k = \frac{\langle g^{k+1}, g^{k+1} - g^k \rangle}{\langle p^k, g^{k+1} - g^k \rangle}$$

Polak-Ribière (PR) Modification:

In the SF formula of β_k , see denominator:

$\langle p^k, g^{k+1} - g^k \rangle = -\langle p^k, g^k \rangle$ (because new gradient is orthogonal to all old conjugate directions). Thus

$\langle p^k, g^{k+1} - g^k \rangle = -\langle -g^k + \beta_{k-1} p^{k-1}, g^k \rangle = \langle g^k, g^k \rangle$ gives:

$$\beta_k = \frac{\langle g^{k+1}, g^{k+1} - g^k \rangle}{\langle g^k, g^k \rangle}$$

Fletcher - Reeves (FR) Modification:

Use the fact that all new gradients are orthogonal to the old ones in the numerator of the PR modification to get:

$$\beta_k = \frac{\langle g^{k+1}, g^{k+1} \rangle}{\langle g^k, g^k \rangle}$$

Important Remark: Stopping criterion.

For nonquadratic problems, the algorithm will not usually converge in n steps, and as the algorithm, the conjugate direction will tend to deteriorate. Thus a common practice is to reinitialize the direction vector to the negative of gradient after every few iteration (e.g. n).