# 1   General

line search conditions:

$$f_{k+1}^T \leq f_k + c_1 \alpha_k \nabla f_k^T p_k, c_1, \alpha_k \in (0,1) \quad (1)$$
$$\nabla f_{k+1}^T p_k \geq c_2 \nabla f_k^T p_k, 0 < c_1 < c_2 < 1 \quad (2)$$
$$\text{where :} \quad (3)$$
$$f_k = f(x_k) \quad (4)$$
$$f_{k+1} = f(x_k + \alpha_k p_k) \quad (5)$$

# 2   Quasi Newton

## 2.1   BFGS

properties: $O(n^2)$, self correcting, slightly more iterations than Newton Method, linear convergence order and superlinear rate of convergence

secant equation:

$$B_{k+1}(x_{k+1} - x_k) = \nabla f_{k+1} - \nabla f_k$$
$$B_{k+1} s_k = y_k$$
$$s_k = \alpha_k p_k$$
$$y_k = \nabla f_{k+1} - \nabla f_k$$

$$B_{k+1} \succ 0$$
$$s_k^T B_{k+1} s_k = s^T y_k > 0$$

*Proof.*

$$y_k^T s_k = (\nabla f_{k+1} - \nabla f_k)^T s_k$$
$$\nabla f_{k+1}^T s_k \geq c_2 \nabla f_k^T s_k$$
$$(\nabla f_{k+1} - \nabla f_k)^T s_k \geq c_2 \nabla f_k^T s_k - \nabla f_k^T s_k$$
$$y_k^T s_k \geq (c_2 - 1)\nabla f_k^T s_k$$
$$c_2 < 1, s_k \text{ is a descent dir} \implies s_k^T y_k > 0$$

Curvature condition holds. $\qquad\square$

constrain $B$ by solving:

$$\min_B ||B - B_k||$$
$$s.t. \ B = B^T, Bs_k = y_k$$

similarly, constrain $B$'s inverse, $H$ where it satisfy secant equation:

$$H_{k+1} y_k = s_k$$
$$\min_H ||H - H_k||$$
$$s.t. \ H = H^T, Hy_k = s_k$$

using weighted Frobenius norm:

$$||A||_W := ||W^{1/2} A W^{1/2}||_F$$
$$||X||_F := (\sum_i \sum_j (X_{ij})^2)^{1/2}$$

solved weight matrix $W$ satisfy $Ws_k = y_k$ solution given by:

$$H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T)$$
$$+ \rho_k s_k^T s_k$$
$$\rho_k = \frac{1}{y_k^T s_k}$$

$W$ is the average Hessian $\bar{G}$:

$$\bar{G} = \int_0^1 \nabla^2 f(x_k + \tau \alpha_k p_k) d\tau$$

initial $H_0$ can be chosen approximately (eg: finite differences, $I$)

---

**Algorithm 1:** BFGS Algorithm

$H_0, x_0, \epsilon > 0$**:** inverse Hessian approx., initial point, convergence tolerance

$x$         **:** solution

1  $k \leftarrow 0$
2  **while** $||\nabla f_k|| > \epsilon$ **do**
3  $\quad \alpha_k \leftarrow \text{LineSearch}(..)$
4  $\quad x_{k+1} \leftarrow x_k + \alpha_k p_k$
5  $\quad s_k \leftarrow x_{x+1} - x_k$
6  $\quad y_k \leftarrow \nabla f_{k+1} - \nabla f_k$
7  $\quad \rho_k \leftarrow \frac{1}{y_k^T s_k}$
8  $\quad H_{k+1} \leftarrow (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T)$
9  $\quad + \rho_k s_k^T s_k$
10  $\quad k \leftarrow k + 1$

---

using Sherman-Morrison-Woodbury formula to obtain Hessian update equation:

$$B_{k+1} = B_k - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \frac{y_k y_k^T}{y_k^T s_k}$$

proper line search is required so that BFGS algo captures curvature information

inaccurate line search can be used to reduce computation cost

# 3    Trust Region Methods

idea:

- models local behaviour of the objective function (eg: 2nd order Taylor series)

- set local region to explore, then simultaneously find direction and step size to take

- region size adaptively set using results from previous iterations

- step may fail due to inadequately set region, which need to be adjusted

- superlinear convergence when approximate model Hessian is equal to true Hessian

using 2nd order Taylor series model with symmetric matrix approximating Hessian

$$\min_{p\in\mathbb{R}^n} m_k(p) = f_k + g_k^T p + \frac{1}{2}p^T B_k p, \ st. \ ||p||\leq \Delta_k$$

$$\Delta_k := \text{trust region radiius}$$
$$g_k = \nabla f(x_k)$$

full step is $(p_k = -B_k^{-1}g_k)$ taken when $B \succ 0$ and $||B_k^{-1}g_k||\leq \Delta_k$

evaluate goodness of model with actual function by:

$$\rho_k = \frac{f(x_k) - f(x_k + p_k)}{m_k(0) - m_k(p_k)}$$

$$action \leftarrow \begin{cases} \text{expand trust region} & , \rho_k \approx 1 \ (agreement) \\ \text{shrink trust region} & , \rho_k < 0 + \text{thresh} \\ \text{keep trust region} & , o/w \end{cases}$$

---
**Algorithm 2:** Trust Region Algorithm

---
**1** $k \leftarrow 0$
**2 while** $||\nabla f_k|| > \epsilon$ **do**
**3** $\quad$ $p_k \leftarrow$
$\qquad$ $\underset{p}{\text{argmin}} f_k + g_k^T p + \frac{1}{2}p^T B_k p, \ st. \ ||p||\leq \Delta_k$
**4** $\quad$ $\rho_k = \frac{f(x_k)-f(x_k+p_k)}{m_k(0)-m_k(p_k)}$
**5** $\quad$ **if** $\rho_k < \gamma(: \frac{1}{4})$ **then**
**6** $\qquad$ $\Delta_{k+1} \leftarrow \alpha(: \frac{1}{4})\Delta_k$
**7** $\quad$ **else if** $\rho_k > \beta(: \frac{3}{4})$ *and* $||p_k||== \Delta_k$ **then**
**8** $\qquad$ $\Delta_{k+1} \leftarrow min(2\Delta_k, \hat{\Delta})$
**9** $\quad$ **else**
**10** $\qquad$ $\Delta_{k+1} \leftarrow \Delta_k$
**11** $\quad$ **if** $\rho_k > \eta(:\in [0, \frac{1}{4}))$ **then**
**12** $\qquad$ $x_{k+1} \leftarrow x_k + p_k$
**13** $\quad$ **else**
**14** $\qquad$ $x_{k+1} \leftarrow x_k$

---

minimizer of the 2nd order Taylor series satisfy the following:

$$(B + \lambda I)p^* = -g$$
$$\lambda(\Delta - ||p^*||) = 0$$
$$(B + \lambda I) \succeq 0$$

solving 2nd order Taylor series using approx methods:

- dogleg

- 2-D subspace minimization

- conjugate gradient based

Cauchy Point:
Use 1st order approx. of model and gradient descent to get get next iterate, bounded within trust region.

## 3.1    Dogleg method

if $B \succ 0$:
$p^B = -B^{-1}g$
$p^* = p^B$ if $\Delta \geq \|p^B\|$
$p^U = \frac{-g^T g}{g^T B g}g$ (intermediate point along direction of steepest descent)
interpolate between $p^U$ and $p^B$:
$$\tilde{p}(\tau) = \begin{cases} \tau p^U & , \tau \in [0, 1] \\ p^U + (\tau - 1)(p^B - p^U) & , \tau \in [1, 2] \end{cases}$$
$B \succ 0 \implies \|\tilde{p}(\tau)\|$ increases wrt. $\tau$, $m(\tilde{p}(\tau))$ decreases wrt. $\tau$

if $\|p^B\| \leq \Delta$: $p$ chosen at $p^B$
else $p$ chosen at intersection of $\tilde{p}(\tau)$ and trust region boundnary by solving:
$\|p^U + (\tau - 1)(p^B - p^U)\|^2 = \|\Delta^2\|$

$p_k^S = \underset{p}{\text{argmin}} f_k + g_k^T p, \|p\| \leq \Delta_k$

$\tau_k = \underset{\tau \geq 0}{\text{argmin}} \, m_k(\tau p_k^S), \|\tau p_k^S\| \leq \Delta_k$

$p_k^S = \frac{-\Delta_k g_k|}{\|g_k\|}$
$p_k^C = \tau_k p_k^S$
$p_k^C = -\tau_k \frac{g_k}{\|g_k\|}$

$$\tau_k = \begin{cases} 1 & , g_k^T B_k g_k \leq 0 \\ min(\frac{\|g_k\|\hat{3}}{\Delta_k g_k^T B_k g_k}, 1) & , o/w \end{cases}$$

## 3.2   Iterative Solution

Idea: solve subproblem $\min_{\|p\| \leq \Delta} m(p)$ by applying Newton's method to find $\lambda$ that matches trust region radius. This is slightly more accurate per step compared to Dogleg. Use $(B + \lambda I)p^* = -g$ to solve $\min_{\|p\| \leq \Delta} m(p)$ for $\lambda$.

If $lambda = 0$ and $(B + \lambda I)p^* = -g, \|p^*\| \leq \Delta$ and $(B + \lambda I) \succeq 0$: return $\lambda$
Else: find $\lambda$ s.t. $(B + I) \succeq 0$ and $\|p(\lambda)\| = \Delta, p(\lambda) = -(B + \lambda I)^{-1}g$. Solve and return $\lambda$.

Solve $\|p(\lambda)\| - \Delta = 0, \lambda > \lambda_1$ via Newton's method (root finding). Approx. this to nearly a linear problem for easy solving:

---

**Algorithm 3:** Subproblem Algo

---

1  **for** $l = 0, 1, ..$ **do**
2  $\quad$ solve $B + \lambda^l I = R^T R$
3  $\quad$ $R^T R p_l = -g$
4  $\quad$ $R^T q_l = p_l$
5  $\quad$ $\lambda^{l+1} \leftarrow \lambda^l + \left(\frac{\|p_l\|}{\|q_l\|}\right)^2 \left(\frac{\|p_l\| - \Delta)}{\Delta}\right)$ check $\lambda \geq \lambda_1$

---

# 4 Conjugate Gradient

## 4.1 linear method

Assuming unconstrained problem with strict convex quadratic objective function:

$$\frac{1}{2}x^T A x - b^T x, A \succ 0, A^T = A$$

$\nabla(\frac{1}{2}x^T A x - b^T x) = Ax - b$, thus $\min_x x^T A x - b^T x$ transformed to $Ax - b = 0$.

$x_{k+1} = x_k + \alpha_k p_k$, solve for $\alpha$:

$$\frac{\partial}{\partial \alpha}((x_k + \alpha_k p_k)^T A(x_k + \alpha_k p_k) - b^T(x_k + \alpha_k p_k)) = 0$$

$$\alpha_k = \frac{-p_k^T r_k}{p_k^T A p_k}$$

## 4.2 Conjugate Direction

Search directions linearly independent wrt. A.

$$(\forall i \neq j) p_i^T A p_j = 0$$

Properties:

- Residual elimnated one direction at a time, resulting in max of n iterations.

- Optimal if Hessian is diagonal, if not can try preconditioning.

- Current residual is orthogonal to all previous search directions.

- Any set of conjugate directions can be used (eg: eigenvectors, Gram-Schmidt)

Expanding subspace minimizer:
Using conjugate directions to generate sequence $\{x\}$, then:
$r_k^T p_i = 0, \forall i < k$, $x_k$ is minimizer of $\frac{1}{2}x^T A x - b^T x$ over $\{x | x = x_0 + span\{p_0, ...p_{k-1}\}$

*Proof.*

$$\tilde{x} = x_0 + \sum_i \sigma_i p_i$$

$\tilde{x}$ minimizes over $\{x_0 + span\{p_0, ...p_{k-1}\}\} \iff r(\tilde{x})^T p_i = 0$

$$h(\sigma) = \phi(\tilde{x})$$

$$\phi(x) = \frac{1}{2}x^T A x - b^T x$$

h is also strictly convex quadratic,
with unique $\sigma^*$ satisfying:

$$\frac{\partial h(\sigma^*)}{\partial \sigma_i} = 0, i = [0, k-1]$$

$$\frac{\partial h(\sigma^*)}{\partial \sigma_i} = \nabla\phi(\tilde{x})^T p_i = 0, i = [0, k-1]$$

$$\nabla\phi(x) = Ax - b = r$$

$$r(\tilde{x})^T p_i = 0, i = [0, k-1]$$

$\square$

$p_i^T r_k = 0, i = [0, k-1]$ via induction:

*Proof.*

base case : $x_1 = x_0 + \alpha_0 p_0$ minimizes $\phi$ along $p_0$

$\implies r_1^T p_0 = 0$

case: $r_{k-1}^T p_i = 0, i = [0, k-2]$ :

$$r_k = r_{k-1} + \alpha_{k-1} A p_{k-1}$$

$$p_{k-1}^T r_k = p_{k-1}^T r_{k-1} + \alpha_{k-1} p_{k-1}^T A p_{k-1} = 0 \text{(by construction)}$$

A-conjugate $\implies p_{k-1}^T A p_{k-1}$

case: $\forall i = [0, k-2] : p_i^T r_k = 0$

$$p_i^T r_k = p_i^T r_{k-1} + \alpha_{k-1} p_i^T A p_{k-1}$$

$$p_i^T r_{k-1} = 0 (by induction hypothesis$$

$$\alpha_{k-1} p_i^T A p_{k-1} = 0 (conjugacy)$$

$$p_i^T r_k = 0, i = [0, k-1]$$

$\square$

## 4.3   Conjugate Gradient Method

Idea:

- uses only previous search direction to compute current search direction

- $p_k$ set to linear combination of $-r_k$ and $p_{k-1}$

- impose $p_k^T A p_{k-1} = 0$

$$p_k = -r_k + \beta_k p_{k-1}$$
$$p_{k-1}^T A p_k = -p_{k-1}^T A r_k + \beta p_{k-1}^T A p_{k-1}$$
$$0 = -p_{k-1}^T A r_k + \beta p_{k-1}^T A p_{k-1}$$
$$\beta = \frac{p_{k-1}^T A r_k}{p_{k-1}^T A p_{k-1}}$$
$$p_0 = -(A x_0 - b) = -r_0$$

---

**Algorithm 4:** Basic Conjugate Gradient Algorithm

---

**1** $r_0 \leftarrow A x_0 - b$
**2** $p_0 = -r_0$
**3** **for** $k = [0, .. n - 1]$ **do**
**4**     **if** $r_k == 0$ **then**
**5**        return $x_k$
**6**     **else**
**7**        $\alpha_k \leftarrow \frac{-r_k^T p_k}{p_x^T A p_k}$
**8**        $x_{k+1} \leftarrow x_k + \alpha_k p_k$
**9**        $r_{k+1} \leftarrow A x_{k+1} - b$
**10**       $b_{k+1} \leftarrow \frac{p_k^T A r_{k+1}}{p_k^T A p_k}$
**11**       $p_{k+1} \leftarrow -r_{k+1} + \beta_{k+1} p_k$

---

$p$ and $r_k$ is within krylov subspace:
$K(r_0; k) = span\{r_0, A r_0, .. A^k r_0\}$
if $r_k \neq 0$:
$r_k^T r_i = 0, i = [0, k - 1]$
$span\{r_0, .., r_k\} = span\{r_0, A r_0, .., A^k r_0\}$
$span\{p_0, .., p + k\} = span\{r_0, A r_0, .., A^k r_0\}$
$p_k^T A p_i = 0, i = [0, k - 1]$
then, $\{x_k\} \rightarrow x^*$ in at most n steps.

Simplification:

$$p_{k+1} \leftarrow -r_{k+1} + \beta_{k+1} p_k$$
$$\alpha_k \leftarrow \frac{-r_k^T p_k}{p_k^T A p_k}$$
$$\alpha_k \leftarrow \frac{-r_k^T (-r_k + \beta_k p_{k-1})}{p_k^T A p_k}$$
$$(\forall i = [0, k-1]) r_k^T p_i = 0 \implies \beta_k r_k^T p_{k-1} = 0$$
$$\alpha_k \leftarrow \frac{r_k^T r_k}{p_k^T A p_k}$$

$$r_{k+1} = r_k + \alpha_k A p_k$$
$$A p_k = \frac{r_{k+1} - r_k}{\alpha_k}$$
$$\beta = \frac{p_k^T A r_{k+1}}{p_k^T A p_k}$$
$$p_k^T A p_k = p_k^T \frac{r_{k+1} - r_k}{\alpha_k} = \frac{-p_k^T r_k}{\alpha} \text{(conjugacy)}$$
$$p_k^T A p_k = -\frac{(-r_k + \beta_k p_{k-1})^T r_k}{\alpha} = \frac{r_k^T r_k}{\alpha} \text{(conjugacy)}$$
$$p_k^T A r_{k+1} = r_{k+1}^T A p_k$$
$$p_k^T A r_{k+1} = r_{k+1}^T \frac{r_{k+1} - r_k}{\alpha_k}$$
$$r_k \in span\{p_k, p_{k-1}\} \text{ and } r_{k+1}^T p_i = 0, i = [0, k] \implies$$
$$p_k^T A r_{k+1} = \frac{r_{k+1}^T r_{k+1}}{\alpha_k}$$
$$\beta_{k+1} \leftarrow \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$$

## 4.4   Nonlinear Method

Minimize general convex function or nonlinear function. Variants: FR, PR.

### 4.4.1   FR

Modify linear CG by:

- replace residual by gradient of objective, $\nabla f_k$

- replace $\alpha_k$ computation by a linear search to find approx. minimum along search direction

Equivalent to linear CG if objective is strongly convex quadratic.

Linear search for $\alpha_k$ with stong Wolfe condition to ensure $p_k$'s are descent directions wrt. objective function.

### 4.4.2 PR

Replace $\beta_{k+1}$ computation in FR with:

$$\beta_{k+1} \leftarrow \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{\nabla f_k^T \nabla f_k}$$