

# 1 General

line search conditions (1st and 2nd order), how to choose step length and direction:

sufficient decrease, proportional to  $\alpha_k$  step length and directional derivative  $\nabla f_k^T p_k$ :

$$f_{k+1}^T \leq f_k + c_1 \alpha_k \nabla f_k^T p_k \quad (1)$$

$$\nabla f_k^T p_k \leq 0 \text{ //by construction} \quad (2)$$

curvature:

$$\nabla f_{k+1}^T p_k \geq c_2 \nabla f_k^T p_k \quad (3)$$

$$\frac{\partial \phi(\alpha_k)}{\partial \alpha_k} \geq \frac{\partial \phi(0)}{\partial \alpha_k} \quad (4)$$

$$c_1, \alpha_k \in (0, 1) \quad (5)$$

$$0 < c_1 < c_2 < 1 \quad (6)$$

where  $f_k = f(x_k)$

$$f_{k+1} = f(x_k + \alpha_k p_k)$$

$$\phi(\alpha_k) = f(x_k + \alpha_k p_k)$$

$$\frac{\partial \phi(\alpha_k)}{\partial \alpha_k} = \nabla f_{k+1}^T p_k$$

$$\frac{\partial \phi(0)}{\partial \alpha_k} = \nabla f_k^T p_k \text{ //initial slope}$$

strong Wolfe curvature condition, to restrict large positive derivative:

$$|\nabla f_{k+1}^T p_k| \leq c_2 |\nabla f_k^T p_k| \quad (7)$$

Form of search direction:

$$p_k = -B_k^{-1} \nabla f_k \quad (8)$$

$B_k$  symmetric, non-singular, postive definite  $\implies p_k$  is a descent direction:

$$\nabla f_k^T (-B_k^{-1} \nabla f_k) < 0 \quad (9)$$

Goldstein condition (may miss minimizer of  $f$ ):

$$f_k + (1 - c) \alpha_k \nabla f_k^T p_k \leq f_{k+1} \leq f_k + c \alpha_k \nabla f_k^T p_k \quad (10)$$

$$c \in (0, \frac{1}{2}) \quad (11)$$

---

## Algorithm 1: Line Search

---

$f, x, d, c_1, \alpha, \beta$ : function, x, direction, gradient threshold, initial step length, contraction

$\alpha$  : found step length

**1 while**  $f(x + \alpha d) > f(x) + c_1 \alpha \nabla f(x)^T d$  **do**

**2**     $\alpha \leftarrow \alpha * \beta$

**3 return**  $\alpha$

---

## 2 Quasi Newton

### 2.1 Concept

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

Derivation, using 2nd order model, with  $B_k$  SPSD:

$$m_k(p) = f_k + \nabla f_k^T p + \frac{1}{2} p^T B_k p$$

taking gradient wrt.  $p$  and solve, assuming minimum exists:

$$0 = \nabla f_k + B_k p$$

$$p = -B_k^{-1} \nabla f_k$$

update equation:

$$x_{k+1} = x_k + \alpha_k p_k$$

$$x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f_k$$

pick  $\alpha$  to satisfy Wolfe conditions

updated model at next iterate  $x_{k+1}$ :

$$m_{k+1}(p) = f_{k+1} + \nabla f_{k+1}^T p + \frac{1}{2} p^T B_{k+1} p$$

impose reasonable conditions:

1. gradient of  $m_{k+1}$  matches gradient of objective function  $f$  at iterate  $x_k$
2. gradient of  $m_{k+1}$  matches gradient of objective function  $f$  at iterate  $x_{k+1}$

$$\nabla m_{k+1}(0) = \nabla f_{k+1} \implies \text{condition 2 satisfied}$$

for condition 1 (gradient of  $m_{k+1}$  match gradient of objective  $f$  at iterate  $x_k$ ):

$$\nabla m_{k+1}(-\alpha_k p_k) = \nabla f_{k+1} - \alpha_k B_{k+1} p_k = \nabla f_k$$

$$\nabla f_{k+1} - \nabla f_k = \alpha_k B_{k+1} p_k$$

$$y_k = \nabla f_{k+1} - \nabla f_k$$

$$s_k = x_{k+1} - x_k = \alpha_k p_k$$

$$y_k = B_{k+1} s_k \text{ [secant equation]}$$

$B_{k+1}$  SPD  $\implies s_k^T B_{k+1} s_k = s_k^T y_k > 0$   
 $f$  strongly convex  $\implies s_k^T y_k > 0$  satisfied for any 2 points  $x_k$  and  $x_{k+1}$ .  
 $f$  nonconvex  $\implies$  need to enforce  $s_k^T y_k$  by Wolfe conditions, and use line search for step length  $\alpha$ .

Using line search ensures curvature condition:

*Proof.*

$$s_k = x_{k+1} - x_k = \alpha_k p_k$$

$$y_k = \nabla f_{k+1} - \nabla f_k$$

Wolfe curvature condition:

$$\nabla f_{k+1}^T p_k \geq c_2 \nabla f_k^T p_k, c_2 \in (0, 1)$$

$$y_k^T s_k = (\nabla f_{k+1} - \nabla f_k)^T (\alpha_k p_k)$$

$$y_k^T s_k = \alpha_k (\nabla f_{k+1} - \nabla f_k)^T p_k$$

$$y_k^T s_k = \alpha_k (c_2 \nabla f_k - \nabla f_k)^T p_k$$

$$y_k^T s_k = \alpha_k (c_2 - 1) \nabla f_k^T p_k$$

$$\alpha > 0 \wedge c_2 - 1 < 0 \wedge \nabla f_k^T p_k < 0 (p_k \text{ is a descent dir})$$

$$\implies y_k^T s_k > 0$$

Curvature condition holds when Wolfe line search is performed.  $\square$

Choosing approximate Hessian,  $B$ :

Force a unique solution among infinite many due to extra degrees of freedom in the matrix, few (n) constraining conditions imposed by second equation, few (n) constraining conditions of PD. One approach is solving a optimization problem to make row rank update to previous iterate:

$$\min_B \|B - B_k\|$$

$$s.t. B = B^T$$

$$B s_k = y_k \text{ [secant equation]}$$

$$B \succ 0$$

alternatively, constrain  $B$ 's inverse,  $H$ :

$$\min_H \|H - H_k\|$$

$$s.t. H = H^T$$

$$H y_k = s_k \quad H_{k+1} y_k = s_k \text{ [secant equation]}$$

Different norms can be used. One choice: weighted Frobenius norm:

$$\|A\|_W := \|W^{1/2} A W^{1/2}\|_F$$

$$\|X\|_F := \left( \sum_i \sum_j (X_{ij})^2 \right)^{1/2}$$

## 2.2 DFT update algorithm

let  $\bar{G}_k$  be the average Hessian:

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

using Taylor's theorem:

$$y_k = \bar{G}_k s_k = \bar{G}_k \alpha_k p_k$$

$$\text{let } W = \bar{G}_k^{-1}$$

solve optimization problem:

$$\min_B \|B - B_k\|_W$$

$$s.t. W = \bar{G}_k^{-1}$$

$$B s_k = y_k$$

$$B = B^T$$

Solution:

$$B_{k+1} = (I - \rho_k y_k s_k^T) B_k (I - \rho_k s_k y_k^T) + \rho_k y_k y_k^T$$

$$\rho_k = \frac{1}{y_k^T s_k}$$

Computation simplification: inverse Hessian used in update of search direction:  $p + k = -B_k^{-1} \nabla f_k$

Use Sherman-Morrison-Woodbury formula to obtain inverse.

let  $H_k = B_k^{-1}$ , then DFP update becomes:

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

This is a rank 2 modification to previous iterate for efficiency.

## 2.3 BFGS update algorithm

Idea: impose condition on inverse of Hessian instead of Hessian approximation.

let  $H$  be approximate inverse of Hessian, then impose:

$$H_{k+1} \succ 0$$

$$H_{k+1} = H_{k+1}^T$$

$$H_{k+1} y_k = s_k$$

solve for  $H$  in the optimization problem:

$$\min_H \|H - H_k\|_W$$

$$s.t. H = H^T$$

$$H y_k = s_k$$

where  $y_k = W s_k$

let  $W$  be the average Hessian:

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

solve to obtain:

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

$$\rho_k = \frac{1}{y_k^T s_k}$$

initial  $H_0$  can be chosen approximately:

- finite differences
- $\alpha I$

---

**Algorithm 2:** BFGS Algorithm

---

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

```

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

---

Cost:  $O(n^2)$  + cost of eval  $f(\cdot)$  + cost of eval  $\nabla f(\cdot)$ .  
Order of convergence: superlinear, worse than Newton but more computationally efficient than Newton.

Using Sherman-Morrison-Woodbury formula to obtain Hessian update equation,

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

, but is it more efficient to use the inverse version.

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

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

$\Delta_k$  := trust region radius

$$g_k = \nabla f(x_k)$$

$$B_k \succeq 0$$

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:

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

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

---

#### Algorithm 3: Trust Region Algorithm

---

```

1  $k \leftarrow 0$ 
2 while  $\|\nabla f_k\| > \epsilon$  do
3    $p_k \leftarrow$ 
      $\underset{p}{\operatorname{argmin}} f_k + g_k^T p + \frac{1}{2} p^T B_k p, \text{ st. } \|p\| \leq \Delta_k$ 
4    $\rho_k = \frac{f(x_k) - f(x_k + p_k)}{m_k(0) - m_k(p_k)}$  //ratio test
5   //change trust region for next iterate
6   if  $\rho_k < \gamma(\frac{1}{4})$  then
7      $\Delta_{k+1} \leftarrow \alpha(\frac{1}{4})\Delta_k$  //shrink
8   else if  $\rho_k > \beta(\frac{3}{4})$  and  $\|p_k\| = \Delta_k$  then
9      $\Delta_{k+1} \leftarrow \min(2\Delta_k, \hat{\Delta})$  //expand
10  else
11     $\Delta_{k+1} \leftarrow \Delta_k$  //keep
12  //accept/reject for current iterate
13  if  $\rho_k > \eta(\in [0, \frac{1}{4}])$  then
14     $x_{k+1} \leftarrow x_k + p_k$ 
15  else
16     $x_{k+1} \leftarrow x_k$ 

```

---

Minimizer of the 2nd order Taylor series satisfy the following as a global solution for trust region iff:

$$(B + \lambda I)p^* = -g$$

complementary slackness:

$$\lambda(\Delta - \|p^*\|) = 0$$

$$(B + \lambda I) \succeq 0$$

$$\lambda \geq 0$$

$\lambda \geq 0 \implies p$  parallel to negative gradient of model

Solving 2nd order Taylor series using approx methods:

- dogleg
- 2-D subspace minimization
- conjugate gradient based, effective when  $B$  large, sparse

These reduce at least as much as Cauchy Point.

Cauchy Point (sufficient reduction in model value):  
Use 1st order approx. of model and gradient descent to get next iterate, bounded within trust region.

#### 3.1 Dogleg method

$$p^B = -B^{-1}g$$

$p^* = p^B$  if  $\Delta \geq \|p^B\|$

$p^U = \frac{-g^T g}{g^T B g} g$  (intermediate point along direction of steepest descent)

Resulting algo, interpolate between  $p^U$  and  $p^B$ :

$$\tilde{p}(\tau) = \begin{cases} \tau p^U & , \tau \in [0, 1] \\ p^U + (\tau - 1)(p^B - p^U) & , \tau \in [1, 2] \end{cases}$$

$$B \succ 0 \implies \|\tilde{p}(\tau)\| \text{ increases wrt. } \tau \wedge m(\tilde{p}(\tau)) \text{ decreases wrt. } \tau$$

if  $\|p^B\| \leq \Delta$ :  $p$  chosen at  $p^B$

else  $p$  chosen at intersection of  $\tilde{p}(\tau)$  and trust region boundary by solving:

$$\|p^U + (\tau - 1)(p^B - p^U)\|^2 = \|\Delta\|^2$$

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

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

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

$$p_k^C = \tau_k p_k^S$$

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

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

### 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.  $(B + \lambda I) \succeq 0$  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 4: Subproblem Algo

---

```

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

---

### 3.3 Trust Region Newton CG Method

Idea: use trust region algo for the outer iteration, use iterative CG based algorithm to solve for inner optimization problem.

Outer iteration: Algo 3.

Inner optimization problem:

$$\begin{aligned} \min_{p \in \mathbb{R}^n} m(p) &= f + g^T p + \frac{1}{2} p^T B p \\ \text{s.t. } \|p\| &\leq \Delta \end{aligned}$$

---

#### Algorithm 5: Trust Region Newton-CG Subproblem (CG Steihaug)

---

```

1  $\epsilon_k = \eta_k \|\nabla f_k\|$ 
2  $z_0 = 0$ 
3  $r_0 = \nabla f_k$ 
4  $d_0 = -r_0 = \nabla f_k$ 
5 if  $\|r_0\| < \epsilon_k$  then
6   return  $p_k = z_0 = 0$ 
7 for  $j = 0, 1, \dots$  do
8   //dir of nonpositive curvature test
9   if  $d_j^T B_k d_j \leq 0$  then
10    return  $\operatorname{argmin}_{p_k} m_k(p_k = z_j + \tau d_j)$  s.t.
         $\|p_k\| = \Delta_k$ 
11     $\alpha_j = \frac{r_j^T r_j}{d_j^T B_k d_j}$ 
12     $z_{j+1} = z_j + \alpha_j d_j$ 
13    //trust region check
14    if  $\|z_{j+1}\| \geq \Delta_k$  then
15      return  $p_k : p_k = z_j + \tau d_j$  s.t.
         $\|p_k\| = \Delta_k, \tau \geq 0$ 
16     $r_{j+1} = r_j + \alpha_j B_k d_j$ 
17    if  $\|r_{j+1}\| < \epsilon_k$  then
18      return  $p_k = z_{j+1}$ 
19     $B_{j+1} = \frac{r_{j+1}^T r_{j+1}}{r_j^T r_j}$ 
20     $\underline{d}_{j+1} = -r_{j+1} + B_{j+1} d_j$ 
```

---

$\epsilon_j$  can be chosen similarly as in Line Search Newton-CG method, where  $\{\eta_k\}$  is the forcing sequence.

## 4 Conjugate Gradient

### 4.1 linear method

Assuming unconstrained problem with strict convex quadratic objective function:

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

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

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

$$\begin{aligned} Ax_{k+1} - b &= 0 \\ A(x_k + \alpha_k p_k) - b &= 0 \\ \alpha_k A p_k &= b - A x_k \\ r_k &= Ax - b \\ \alpha_k A p_k &= -r_k \\ \alpha_k p_k^T A p_k &= -p_k^T r_k \\ \alpha_k &= -\frac{p_k^T r_k}{p_k^T A p_k} \end{aligned}$$

### 4.2 Conjugate Direction

Enforce by construction for search directions linearly independent wrt.  $A$ :

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

Properties:

- Residual eliminated 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:
  - eigenvectors
  - Gram-Schmidt
  - conjugate gradient algo.

#### 4.2.1 Termination Steps

**Theorem 4.1** (Conjugate Direction Termination). *Conjugate direction algorithm converges to solution  $x \in \mathbb{R}^n$  of linear system in  $n$  steps. Theorem [T5.1] in Num. Opt. book.*

*Proof.*

given:

$$x_{k+1} = x_k + \alpha_k p_k$$

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

$$(\forall i \neq j) p_i^T A p_j = 0 \text{ [A conjugate]}$$

then:

$$x^* - x_0 = \sum_{i=0}^{n-1} \sigma_i p_i$$

$$p_k^T A(x^* - x_0) = p_k^T A\left(\sum_{i=0}^{n-1} \sigma_i p_i\right)$$

$$\sigma_i = \frac{p_k^T A(x^* - x_0)}{p_k^T A p_k} \text{ (from conjugacy)}$$

$$x_k = x_0 + \sum_{i=0}^{k-1} \alpha_i p_i$$

$$p_k^T A x_k = p_k^T A(x_0 + \alpha_0 p_0 + \dots)$$

$$p_k^T A x_k = p_k^T A x_0 \text{ (from conjugacy)}$$

$$p_k^T A(x_k - x_0) = 0$$

$$p_k^T A(x^* - x_0) = p_k^T A(x^* - x_k)$$

$$p_k^T A(x^* - x_0) = p_k^T (b - A x_k) = -p_k^T r_k$$

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

$\alpha_k$  is 1D minimizer of  $k$ th coordinate by construction

Conjugate direction algo. terminates in  $n$  steps.  $\square$

If  $A$  is not diagonal, can transform coordinates:

$$S \hat{x} = x$$

$$S = [p_0 \ p_1 \ \dots \ p_{n-1}], (\forall i \neq j) p_i^T A p_j = 0$$

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

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

$S^T A S$  diagonal  $\implies$  can optimize one coordinate at a time

### 4.2.2 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  $\frac{1}{2}x^T Ax - b^T x$   
over  $\{x | x = x_0 + \text{span}\{p_0, \dots, p_{k-1}\}\}$

*Proof.*

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

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

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

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

$h$  is also strictly convex quadratic,

with unique  $\sigma^*$  satisfying:

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

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

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

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

□

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

*Proof.*

base case of  $k=1, i=0$  :  $x_1 = x_0 + \alpha_0 p_0$  minimizes  $\phi$  along  $p_0$

$$\implies r_1^T p_0 = (r_0 + \alpha_0 A p_0)^T p_0 = 0$$

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

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

$$(\forall i \in [0, k-2]) p_i^T r_k = p_i^T r_{k-1} + \alpha_{k-1} p_i^T A p_{k-1}$$

$$(\forall i \in [0, k-2]) p_i^T A p_{k-1} = 0 \text{ (A-conjugacy by construction)}$$

$$(\forall i \in [0, k-2]) p_i^T r_{k-1} = 0 \text{ (by induction hypothesis)}$$

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

□

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

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

---

**Algorithm 6:** 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      $\alpha_k \leftarrow \frac{-r_k^T p_k}{p_k^T A p_k} = \frac{r_k^T r_k}{p_k^T A p_k}$ 
8      $x_{k+1} \leftarrow x_k + \alpha_k p_k$ 
9      $r_{k+1} \leftarrow A x_{k+1} - b$ 
10     $\beta_{k+1} \leftarrow \frac{p_k^T A r_{k+1}}{p_k^T A p_k} = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$ 
11     $p_{k+1} \leftarrow -r_{k+1} + \beta_{k+1} p_k$ 

```

---

$p_k$  and  $r_k$  is within krylov subspace:

$$K(r_0; k) = \text{span}\{r_0, A r_0, .. A^k r_0\}$$

if  $r_k \neq 0$ :

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

$$\text{span}\{r_0, .., r_k\} = \text{span}\{r_0, A r_0, .., A^k r_0\}$$

$$\text{span}\{p_0, .., p + k\} = \text{span}\{r_0, A r_0, .., A^k r_0\}$$

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

then,  $\{x_k\} \rightarrow x^*$  in at most n steps.

Simplification:

$$p_{k+1} \leftarrow -r_{k+1} + \beta_{k+1} p_k$$

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

$$\alpha_k \leftarrow \frac{-r_k^T (-r_k + \beta_k p_{k-1})}{p_k^T A p_k}$$

$$(\forall i = [0, k - 1]) r_k^T p_i = 0 \implies \beta_k r_k^T p_{k-1} = 0$$

$$\alpha_k \leftarrow \frac{r_k^T r_k}{p_k^T A p_k} \text{ (simplified)}$$

$$r_{k+1} = r_k + \alpha_k A p_k$$

$$A p_k = \frac{r_{k+1} - r_k}{\alpha_k}$$

$$\beta = \frac{p_k^T A r_{k+1}}{p_k^T A p_k}$$

$$p_k^T A p_k = p_k^T \frac{r_{k+1} - r_k}{\alpha_k} = \frac{-p_k^T r_k}{\alpha_k} \text{ (conjugacy)}$$

$$p_k^T A p_k = -\frac{(-r_k + \beta_k p_{k-1})^T r_k}{\alpha_k} = \frac{r_k^T r_k}{\alpha_k} \text{ (conjugacy)}$$

$$p_k^T A r_{k+1} = r_{k+1}^T A p_k$$

$$p_k^T A r_{k+1} = r_{k+1}^T \frac{r_{k+1} - r_k}{\alpha_k}$$

$$r_k \in \text{span}\{p_k, p_{k-1}\} \text{ and } r_{k+1}^T p_i = 0, i = [0, k] \implies$$

$$p_k^T A r_{k+1} = \frac{r_{k+1}^T r_{k+1}}{\alpha_k}$$

$$\beta_{k+1} \leftarrow \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k} \text{ (simplified)}$$

### 4.4 Nonlinear Method

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

#### 4.4.1 FR (Fletcher Reaves)

Modify linear CG by:

- replace residual by gradient of nonlinear objective,  $r_k \rightarrow \nabla f_k$
- replace  $\alpha_k$  computation by a linear search to find approx. minimum along search direction of  $f$

Equivalent to linear CG if objective is strongly convex quadratic.

Linear search for  $\alpha_k$  with strong 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:

$$\begin{aligned}\beta_{k+1}^{PR} &\leftarrow \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{\nabla f_k^T \nabla f_k} \\ \beta_{k+1}^+ &\leftarrow \max(\beta_{k+1}^{PR}, 0)\end{aligned}$$

CG Gradient Algo. Property [T.5.3]:

$$\begin{aligned}(\forall i = [0, k-1]) \quad r_k^T r_i &= 0 \\ \text{span}\{r_0, \dots, r_k\} &= \text{span}\{r_0, Ar_0, \dots, A^k r_0\} \\ \text{span}\{p_0, \dots, p_k\} &= \text{span}\{r_0, Ar_0, \dots, A^k r_0\} \\ (\forall i = [0, k-1]) \quad p_k^T A p_i &= 0 \\ \implies \{x_k\} &\text{converges to } x^* \text{ in at most } n \text{ steps}\end{aligned}$$

Proof by induction to show generated search directions are A-conjugate. Then apply Theorem 4.1 to conclude algo terminates within n steps.

## 5 Proximal Algorithm

Idea:

- reliance on easy to evaluate proximal operators
- separability allows parallel evaluation
- generalization of projection based algorithms

$$\text{prox}_{\lambda f}(v) = \underset{x}{\operatorname{argmin}} f(x) + \frac{1}{2\lambda} \|x - v\|_2^2$$

Resolvent of subdifferential operator:

$$\begin{aligned}z = \text{prox}_{\lambda f}(x) &\implies z \in (I + \lambda \partial f)^{-1}(x) \\ (I + \lambda \partial f)^{-1} &:= \text{resolvent of operator } \partial f\end{aligned}$$

### 5.1 Proximal Gradient Method

Solve  $\min_x g(x) + f(x)$ , where  $f, g$  are closed, convex functions and  $f$  differentiable

$$\begin{aligned}x^* &= \text{prox}_{\lambda^k g}(x^k - \lambda^k \nabla f(x^k)) \\ &= \underset{x}{\operatorname{argmin}} g(x) + \frac{1}{2\lambda^k} \|x - (x^k - \lambda^k \nabla f(x^k))\|_2^2\end{aligned}$$

tradeoff between  $g$  and gradient step

$g = I_C(x) \implies$  projected gradient step

$g = 0 \implies$  gradient descent

$f = 0 \implies$  proximal minimization

Relation to Fixed Point:

$x^*$  is a fixed point solution of  $\min_x g(x) + f(x)$  iff  $0 \in \nabla f(x^*) + \partial g(x^*)$  iff  $x^* = (I + \lambda \partial g)^{-1}(I - \lambda \nabla f)(x^*)$

Forward Euler, Backward Euler stepping is same as the proximal gradient iteration,  $\text{prox}_{\lambda^k g}(x^k - \lambda^k \nabla f(x^k))$

### 5.2 Accelerated Proximal Gradient Method

Introduce extrapolation:

$$\begin{aligned}y^{k+1} &= x^k + w^k(x^k - x^{k-1}) \\ x^{k+1} &= \text{prox}_{\lambda^k g}(y^{k+1} - \lambda^k \nabla f(y^{k+1})) \\ w^k &\in [0, 1)\end{aligned}$$

Example:  $w^k = \frac{k}{k+3}, w^0 = 0, \lambda^k \in (0, 1/L], L :=$  Lipschitz constant of  $\nabla f$ , or  $\lambda^k$  found via line search. Line search for  $\lambda^k$  (Beck and Teboulle):

**Algorithm 7:** Proximal Gradient Algorithm

---

```

1  $\hat{f}(x, y) := f(y) + \nabla f(y)^T(x - y) + \frac{1}{2\lambda}\|x - y\|_2^2$ 
2 while True do
3    $z = \text{prox}_{\lambda g}(y^k - \lambda \nabla f(y^k))$ 
4   if  $f(x) \leq \hat{f}(z, y^k)$  then
5      $\lambda = \beta \lambda$ 
6    $\lambda = \beta \lambda$ 
7 return  $\lambda^k := \lambda, x^{k+1} := z$ 

```

---

**5.3 Types of Proximal Operators**

- quadratic functions

$$f = \frac{1}{2}\|\cdot\|_x^2 \implies \text{prox}_{\lambda f}(v) = \left(\frac{1}{1+\lambda}\right)v$$

$$f = \frac{1}{2}x^T A x + b^T x + c, A \in S_+^n \implies$$

$$\text{prox}_{\lambda f}(v) = (I + \lambda)^{-1}(v - \lambda b)$$

- unconstrained problem: use gradient methods such as Newton, Quasi-Newton
- constrained: use projected subgradient for non-smooth, projected gradient or interior method for smooth
- separable function: if scalar, may be solved analytically, eg: L1 norm separable to:

$$f(x) = |x| \implies \text{prox}_{\lambda f}(v) = \begin{cases} v - \lambda, & v \geq \lambda \\ 0, & |v| \leq \lambda \\ v + \lambda, & v \leq -\lambda \end{cases}$$

$$f(x) = -\log(x) \implies \text{prox}_{\lambda f}(v) = \frac{v + \sqrt{v^2 + 4\lambda}}{2}$$

- general scalar function

- localization: using a subgradient oracle and bisection algorithm
- twice continuously differentiable: guarded Newton method

- polyhedra constraint, quadratic objective: solve as QP problem

- duality to reduce number of variables to solve if possible
- gram matrix caching

- affine constraint( $Ax = b$ ): use pseudo-inverse,  $A^+$ :

$$\Pi_C(v) = v - A^+(Av - b)$$

$$A \in \mathbb{R}^{m \times n}, m < n \implies A^+ = A^T(AA^T)^{-1}$$

$$A \in \mathbb{R}^{m \times n}, m > n \implies A^+ = (A^T A)^{-1} A^T$$

- hyperplane constant( $a^T x = b$ ):

$$\Pi_C(v) = v + \left(\frac{b - a^T v}{\|a\|_2^2}\right)a$$

- halfspace

$$\Pi_C(v) = v - \frac{\max(a^T v - b, 0)}{\|a\|_2^2}a$$

- box( $l \leq x \leq u$ )

$$\Pi_C(v)_k = \min(\max(v_k, l_k), u_k)$$

- probability simplex( $1^T x = 1, x \geq 0$ )  
bisection algo on  $\nu$ :

$$\Pi_C(v) = (v - \nu 1)_+$$

$$\text{initial } [l_k, u_k] = [\max_i v_i - 1, \max_i v_i]$$

analytically solve when bounded in between 2 adjacent v's

- cones ( $\kappa$ : proper cone)  
problem of the form:

$$\min_x \|x - v\|_2^2$$

$$\text{s.t. } : x \in \kappa$$

$$x \in \kappa$$

$$v = x - \lambda$$

$$\lambda \in \kappa^*$$

$$\lambda^T x = 0$$

- cone  $C = \mathbb{R}_+^n$

$$\Pi_C(v) = v_+$$

- 2nd order cone  $C = \{(x, t) \in \mathbb{R}^{n+1} : \|x\|_2 \leq t\}$

$$\Pi_C(v, s) = \begin{cases} 0, & \|v\|_2 \leq -s \\ (v, s), & \|v\|_2 \leq s \\ \frac{1}{2}(1 + \frac{s}{\|v\|_2})(v, \|v\|_2), & \|v\|_2 \geq |s| \end{cases}$$

- PSD cone  $S_+^n$

$$\Pi_C(V) = \sum_i (\lambda_i)_+ u_i u_i^T$$

$$V = \sum_i \lambda_i u_i u_i^T \text{ (eigendecomposition)}$$

- exponential cone
  - Todo
- pointwise supremum
  - max function
  - support function
- norms
  - L2
  - L1
  - L-inf
  - elastic net
  - sum of norms
  - matrix norm
- sublevel set
- epigraph
- matrix functions
- Todo

## 6 Smoothness

$\beta$  smoothness:

$\frac{\beta}{2}\|x\|_2^2 - f(x)$  is convex (fit quadratic on top)

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{\beta}{2}\|y - x\|_2^2 \\ \implies f \text{ is smooth}$$

$\alpha$  strong convexity:

$f(x) - \frac{\alpha}{2}\|x\|_2^2$  is convex (fit quadratic below)

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) + \frac{\alpha}{2}\|y - x\|_2^2 \\ f \text{ may not be smooth}$$

Smoothness near optimal point: there always exists a step (length  $< \frac{2}{\beta}$ ) that is a descent direction.

Property of  $f$  wrt. optimality value:  
 $f$  is  $\beta$ -smooth  $\implies$

$$\frac{1}{2\beta}\|\nabla f(x)\|_2^2 \leq f(x) - f(x^*) \leq \frac{\beta}{2}\|x - x^*\|_2^2$$

$f$  is  $\alpha$ -strongly convex  $\implies$

$$\frac{\alpha}{2}\|x - x^*\|_2^2 \leq f(x) - f(x^*) \leq \frac{1}{2\alpha}\|\nabla f(x)\|_2^2$$

Co-coercivity for  $\beta$ -smooth  $f$ :

$$(\nabla f(x) - \nabla f(y))^T(x - y) \geq \frac{1}{\beta}\|\nabla f(x) - \nabla f(y)\|_2^2$$

Coercivity for  $\alpha$ -strongly convex  $f$ :

$$(\nabla f(x) - \nabla f(y))^T(x - y) \geq \alpha\|x - y\|_2^2$$

### 6.1 Oracle based lower bounds

Lipschitz convex function: error  $= \mathcal{O}(\frac{1}{\sqrt{T}})$

Smooth convex function:  $\epsilon = \mathcal{O}(\frac{1}{T^2})$

Smooth + strongly convex function:  $\epsilon = \mathcal{O}(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1})$

### 6.2 Gradient descent rate of convergence

$\beta$ -smooth  $f$ :  $\mathcal{O}(\frac{1}{T})$

$\beta$ -smooth +  $\alpha$ -strongly convex  $f$ :  $\mathcal{O}((\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1})^T)$

### 6.3 Gradient descent with momentum

$f$   $\beta$ -smooth,  $\alpha$ -s.c.:

$$x_1 = y_1 = x_{init} \\ y_{t+1} = x_t - \frac{1}{\beta}\nabla f(x_t) \\ x_{t+1} = \left(1 + \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)y_{t+1} - \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)y_t$$

$f$   $\beta$ -smooth:

$$x_{t+1} = x_t - \eta\nabla f(x_t) \text{ [normal GD]} \\ \text{momentum, nesterov accel.:} \\ d_{t+1} = \gamma_{t+1}(x_{t+1} - x_t) \\ x_{t+1} = (x_t + d_t) - \eta\nabla f(x_t + d_t) \\ \epsilon = \mathcal{O}\left(\frac{1}{T^2}\right)$$

## 7 Common Algorithms

### 7.1 Projected Gradient

$$y_{t+1} = x_t - \eta g_t, g_t \in \partial f(x_t) \\ x_{t+1} = \text{proj}_X(y_{t+1}) \\ x^* = \text{proj}_X(x^*) \text{ [fixed point at optimality]}$$

### 7.2 Proximal Gradient

$$\text{prox}_{\eta h}(y) = \underset{x}{\operatorname{argmin}} h(x) + \frac{1}{2\eta}\|x - y\|_2^2 \\ x_{t+1} = \text{prox}_{\eta h}(x_t - \eta\nabla f(x_t))$$

Special case of projected gradient:

let  $h(x) = I_X(x) \implies \text{prox}_{\eta h} = \text{proj}_X$

where  $I_X(x) = \begin{cases} 0 & , x \in X \\ +\infty & , o/w \end{cases}$

#### 7.2.1 Common proximal operators

$h(x) = \|x\|_1$ :

$$(\text{prox}_{\eta h}(x))_i = \begin{cases} x_i - \eta & , x_i \geq \eta \\ 0 & , |x_i| \leq \eta \\ x_i + \eta & , x_i \leq -\eta \end{cases} \\ = \max(|x_i| - \eta, 0) \operatorname{sign}(x_i)$$

$$h(x) = \frac{1}{2}x^T Q x + q^T x + q_0, Q \succeq 0:$$

$$\text{prox}_{\eta h} = (I + \eta Q)^{-1}(x - \eta q)$$

$$h(x) = \sum_i h_i(x_i):$$

$$(prox_{\eta h}(x))_i = prox_{\eta h_i}(x_i) \text{ [parallelism]}$$

Composite function:

$\min_x f(x) = g(x) + h(x)$ ,  $g$  smooth,  $h$  not smooth but has prox. operator.

Convergence due to non-smooth function:  $\mathcal{O}(\frac{1}{\sqrt{T}})$

### 7.3 L1 regularization with subgradient descent

$$\min_x \|Ax - y\|_2^2 + \|x\|_1$$

$$g = \|Ax - y\|_2^2$$

$$h = \lambda \|x\|_1$$

$$\partial_x h(x) = \begin{cases} -1 & , x_i < 0 \\ [-1, 1] & , x_i = 0 \\ 1 & , x_i > 0 \end{cases}$$

$$x_{t+1} = x_t - \eta(\nabla g(x_t) + \lambda \partial h(x_t))$$

$$x_{t+1} = x_t - \eta(2A^T(Ax - y) + \lambda z)$$

$$z = \begin{cases} -1 & , x_i < 0 \\ [-1, 1] & , x_i = 0 \\ 1 & , x_i > 0 \end{cases}$$

Convergence:  $\epsilon = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$

### 7.4 ISTA (with proximal gradient)

$$\min_x g(x) + h(x), g \text{ } \beta\text{-smooth}$$

$$g(x) = \|Ax - y\|_2^2$$

$$h(x) = \lambda \|x\|_1$$

$$x_{t+1} = prox_{\eta h}(x_t - \eta \nabla g(x_t))$$

$$g \text{ } \beta\text{-smooth, select } \eta = \frac{1}{\beta}$$

$$x_{t+1} = prox_{\frac{1}{\beta} \lambda \|\cdot\|_1} \left( x_t - \frac{1}{\beta} \nabla g(x_t) \right)$$

$$x_{t+1} = \operatorname{argmin}_x \lambda \|x\|_1 + \frac{\beta}{2} \|x - (x_t - \frac{1}{\beta} \nabla g(x_t))\|_2^2$$

$$\nabla g(x_t) = 2A^T(Ax - y)$$

Convergence:  $\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$

### 7.5 FISTA

$$\min_x g(x) + h(x), g \text{ } \beta\text{-smooth}$$

$$\mu_0 = 0$$

$$\mu_t = \frac{1 + \sqrt{1 + 4\mu_{t-1}^2}}{2}$$

$$\gamma_t = \frac{1 - \mu_t}{\mu_{t+1}}$$

$$x_1 = z_1 = \text{arbitrary}$$

$$x_{t+1} = (1 - \gamma_t)x_{t+1} + \gamma_t z_t$$

$$z_{t+1} = prox_{\frac{\lambda}{\beta} \|\cdot\|_1} \left( x_t - \frac{1}{\beta} \nabla g(x_t) \right)$$

Convergence:  $\epsilon = \mathcal{O}\left(\frac{1}{T^2}\right)$

## 8 Subgradient Method

TODO

## 9 Large Scale Optimization

1. memory reduction techniques with quasi-Newton methods.
2. Inexact Newton methods with superlinear local convergence properties, indefinite Hessian, Hessian-free impl.
3. function that is partial separable to use smaller subspaces.

### 9.1 Inexact Newton Methods

Iterative method with local convergence properties.

Residual:  $r_k = \nabla^2 f_k p_k + \nabla f_k$

Termination condition:  $\|r_k\| \leq \eta_k \|\nabla f_k\|, \eta_k \in (0, 1)$

Convergence properties: see Thm 7.1, 7.2 (Nocedal)

2 such algorithms:

1. Line Search Newton-CG
2. Trust Region Newton-CG (Section 3.3)

#### 9.1.1 Line Search Newton-CG

Augment original CG algorithm with line search, curvature check and forcing function check.