

Solutions and Notes to Selected Problems In:  
Numerical Optimization  
by Jorge Nocedal and Stephen J. Wright.

John L. Weatherwax\*

December 12, 2019

---

\*wax@alum.mit.edu

## Chapter 2 (Fundamentals of Unconstrained Optimization)

### Problem 2.1

For the Rosenbrock function

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2,$$

we have that (recall that the gradient is a column vector)

$$\begin{aligned}\nabla f(x) &= \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 200(x_2 - x_1^2)(-2x_1) + 2(1 - x_1)(-1) \\ 200(x_2 - x_1^2) \end{bmatrix} \\ &= \begin{bmatrix} -400x_1(x_2 - x_1^2) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix}.\end{aligned}$$

Next for the Hessian we compute

$$\begin{aligned}\nabla^2 f(x) &= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} \\ &= \begin{bmatrix} -400(x_2 - x_1^2) - 400x_1(-2x_1) - 2(-1) & -400x_1 \\ & -400x_1 & 200 \end{bmatrix} \\ &= \begin{bmatrix} -400x_2 + 1200x_1^2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix}.\end{aligned}$$

By the first-order necessary conditions  $\nabla f(x^*) = 0$  for  $x^*$  to be a local minimizer. For this to happen from the second equation in the system  $\nabla f(x) = 0$  we must have

$$x_2 = x_1^2.$$

If we put this into the first equation in the system  $\nabla f(x) = 0$  we have

$$-2(1 - x_1) = 0 \quad \text{so} \quad x_1 = 1.$$

Using the first equation derived this means that  $x_2 = 1^2 = 1$ .

Next, evaluating the Hessian at this point gives

$$\nabla^2 f(x^*) = \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix}.$$

This matrix has two positive eigenvalues and is thus positive definite.

### Problem 2.2

For this function I find

$$\nabla f(x) = \begin{bmatrix} 8 + 2x_1 \\ 12 - 4x_2 \end{bmatrix}.$$

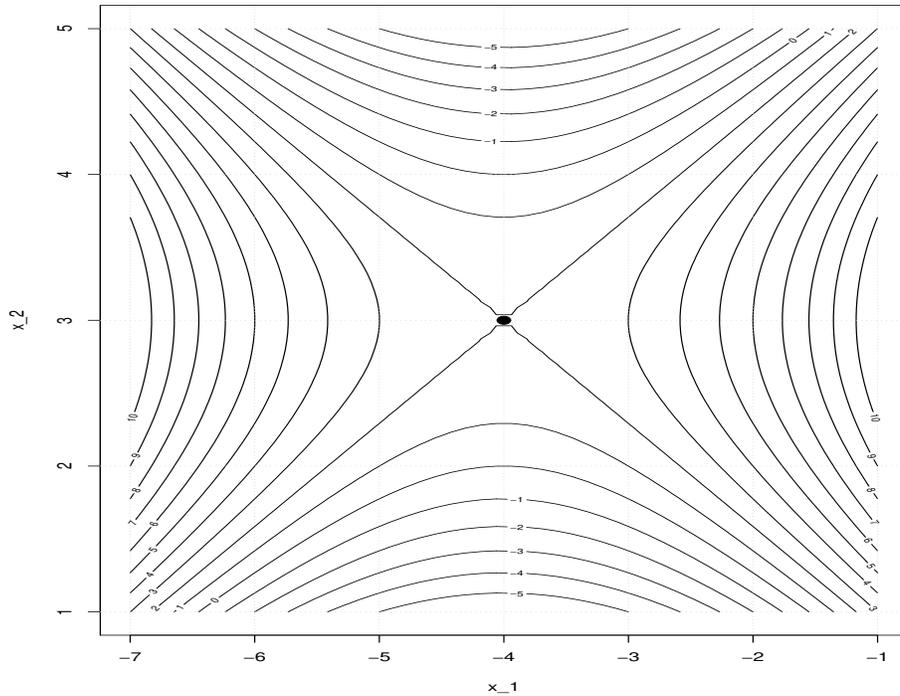


Figure 1: The requested contour plot.

To find the stationary points we set this equal to zero and solve for  $(x_1, x_2)$ . I find  $x_1 = -4$  and  $x_2 = 3$ . From the above form of  $\nabla f(x)$  we have

$$\nabla^2 f = \begin{bmatrix} 2 & 0 \\ 0 & -4 \end{bmatrix}.$$

As this matrix has two eigenvalues of opposite signs every point is a saddle point.

In Figure 1 I present the contour plot for this function centered on the stationary point  $(-4, 3)$ . Looking at the numbers on the contours we see that moving North or South the value of  $f$  decreases while moving East or West the value of  $f$  increases. This is the definition of a “saddle point”.

### Problem 2.3

For  $f_1(x)$  note that  $\nabla f_1(x) = a$  where  $a$  is a  $n \times 1$  vector. From this we have that  $\nabla^2 f_1(x) = 0$  where in this case 0 denotes the  $n \times n$  matrix.

For  $f_2(x)$  note that we can write it as

$$f_2(x) = \sum_{i,j} x_i a_{ij} x_j.$$

From this and using the “Kronecker delta” we have

$$\begin{aligned}
 \frac{\partial f_2(x)}{\partial x_k} &= \sum_{i,j} \delta_{ik} a_{ij} x_j + \sum_{i,j} x_i a_{ij} \delta_{jk} \\
 &= \sum_j a_{kj} x_j + \sum_i x_i a_{ik} \\
 &= \sum_j a_{kj} x_j + \sum_i (a^T)_{ki} x_i \\
 &= (Ax)_k + (A^T x)_k,
 \end{aligned}$$

where the notation  $(\cdot)_k$  means the  $k$ th component of the vector inside the parenthesis. This means that in matrix form we can write

$$\nabla f_2(x) = (A + A^T)x.$$

We now seek to evaluate  $\nabla^2 f_2(x)$ . To compute the  $ij$ th element of that matrix we compute

$$\begin{aligned}
 (\nabla^2 f_2(x))_{ij} &= \frac{\partial}{\partial x_i} ((Ax)_j + (A^T x)_j) \\
 &= \frac{\partial}{\partial x_i} \left( \sum_k a_{jk} x_k + \sum_k a_{kj} x_k \right) \\
 &= \sum_k a_{jk} \delta_{ki} + \sum_k \delta_{ik} a_{kj} \\
 &= a_{ji} + a_{ij}.
 \end{aligned}$$

Thus in matrix form we have

$$\nabla^2 f_2(x) = A + A^T.$$

#### Problem 2.4

For the function  $f(x) = \cos\left(\frac{1}{x}\right)$  we have

$$\begin{aligned}
 f'(x) &= -\sin\left(\frac{1}{x}\right) \left(-\frac{1}{x^2}\right) = \frac{1}{x^2} \sin\left(\frac{1}{x}\right) \\
 f''(x) &= \frac{1}{x^2} \left(\cos\left(\frac{1}{x}\right)\right) \left(-\frac{1}{x^2}\right) - 2 \left(\frac{1}{x^3}\right) \sin\left(\frac{1}{x}\right) \\
 &= -\frac{1}{x^4} \cos\left(\frac{1}{x}\right) - \frac{2}{x^3} \sin\left(\frac{1}{x}\right).
 \end{aligned}$$

For  $x \neq 0$  a second-order Taylor series expansion is thus given by

$$\begin{aligned}
 \cos\left(\frac{1}{x+p}\right) &= \cos\left(\frac{1}{x}\right) + \frac{1}{x^2} \sin\left(\frac{1}{x}\right) p + \frac{1}{2} f''(x+tp) p^2 \\
 &= \cos\left(\frac{1}{x}\right) + \frac{1}{x^2} \sin\left(\frac{1}{x}\right) p - \frac{1}{2} \left( \frac{1}{(x+tp)^4} \cos\left(\frac{1}{x+tp}\right) + \frac{2}{(x+tp)^3} \sin\left(\frac{1}{x+tp}\right) \right),
 \end{aligned}$$

for some  $t \in (0, 1)$ .

Next, for the function  $f(x) = \cos(x)$  we have

$$\begin{aligned}f'(x) &= -\sin(x) \\f''(x) &= -\cos(x) \\f'''(x) &= \sin(x).\end{aligned}$$

Using these a third order Taylor series expansion gives

$$\cos(x + p) = \cos(x) - \sin(x)p - \frac{1}{2} \cos(x)p^2 + \frac{1}{6} \sin(x + tp)p^3,$$

for some  $t \in (0, 1)$ . When  $x = 1$  this becomes

$$\cos(1 + p) = \cos(1) - \sin(1)p - \frac{1}{2} \cos(1)p^2 + \frac{1}{6} \sin(1 + tp)p^3.$$

### Problem 2.5

Our function  $f(x)$  is  $f(x) = \|x\|^2 = x_1^2 + x_2^2$ . Then as  $\cos^2(k) + \sin^2(k) = 1$  we see that

$$f(x_k) = 1 + \frac{1}{2^k}.$$

Then as

$$1 + \frac{1}{2^{k+1}} < 1 + \frac{1}{2^k},$$

we see that  $f(x_{k+1}) < f(x_k)$ .

Now if we are given a point on the unit circle  $\|x\| = 1$  then using this points representation as a polar complex number we have  $x = e^{i\theta}$  for some  $\theta$ . Note that the points in our sequence  $x_k$  can be written as a polar complex number as

$$x_k = \left(1 + \frac{1}{2^k}\right) e^{ik}.$$

From the periodicity of the trigonometric function we have  $e^{ik} = e^{i\xi_k}$  where  $\xi_k$  is given as in the book. Then using the statement in the book (that every value  $\theta$  is a limit point of the sequence  $\xi_k$ ) we have the desired result.

### Problem 2.6

The definition of an isolated local minimizer is given in the book. The fact that  $x^*$  is isolated means that it is the only minimizer in a neighborhood  $\mathcal{N}$  so that  $f(x) > f(x^*)$  for all points  $x \neq x^*$ . This later statement is the fact that  $x^*$  is a strict local minimizer.

### Problem 2.7

Let  $G = \{x_i\}$  be the set of global minimums of our function  $f(x)$ . Then this means that

$$f(x_i) \leq f(x),$$

for all  $x$  and each  $x_i \in G$ . In that relationship take  $x = x_j \in G$  for  $j \neq i$  and conclude that  $f(x_i) \leq f(x_j)$ . We could do the same thing with  $i$  and  $j$  switched to show that  $f(x_j) \leq f(x_i)$ . This means that  $f(x_i) = f(x_j)$  and thus all global minimums must have the same function value. Lets call this value  $m$  so that  $f(x_i) = m$  for all  $x_i \in G$ . Next consider a new point  $z$  from  $x_i$  and  $x_j$  for  $i \neq j$  given by

$$z = \lambda x_i + (1 - \lambda)x_j.$$

As  $f$  is convex we have that

$$f(z) = f(\lambda x_i + (1 - \lambda)x_j) \leq \lambda f(x_i) + (1 - \lambda)f(x_j) = \lambda m + (1 - \lambda)m = m.$$

As  $z$  cannot have  $f(z)$  smaller than the global minimum  $m$  (otherwise  $m$  would not be the true global minimum) we see that  $z$  must be equal to  $m$  and  $z$  is the location of another global minimum thus  $z \in G$ . This shows that the set  $G$  convex set.

### Problem 2.8

To be a decent direction  $p$  at  $x$  means that  $p^T \nabla f(x) < 0$ . For the given function  $f$  we have

$$\nabla f = \begin{bmatrix} 2(x_1 + x_2^2) \\ 2(x_1 + x_2^2)(2x_2) \end{bmatrix}.$$

At the point  $x^T = (1, 0)$  we see that  $\nabla f = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$  and that  $p^T \nabla f = -1(2) + 1(0) = -2 < 0$ .

Thus  $p$  is a decent direction.

The minimizers for the books Eq. 2.9 are to find

$$\min_{\alpha > 0} f(x_k + \alpha p_k).$$

For this problem

$$x + \alpha p = \begin{bmatrix} 1 - \alpha \\ \alpha \end{bmatrix},$$

so that  $f(x + \alpha p) = (1 - \alpha + \alpha^2)^2$ . From this I find

$$\frac{df}{d\alpha} = 2(1 - \alpha + \alpha^2)(-1 + 2\alpha).$$

To find the extremes of this function we need to have the derivative of  $f$  with respect to  $\alpha$  equal to zero which can happen if

$$\alpha = \frac{1}{2},$$

or if the quadratic factor in its representation is zero. The quadratic factor being zero gives complex roots for  $\alpha$  and cannot be zero for  $\alpha > 0$ .

We can show that  $f''(\frac{1}{2}) > 0$  showing that  $\alpha = \frac{1}{2}$  is a minimum of  $f$ .

## Problem 2.9

In the notation of this problem  $\tilde{f}(z)$  means to view the function as a function of the variable  $z$  and the notation  $f(x)$  means to view our function as a function of the variable  $x$ . Of course  $\tilde{f}(z) = f(x)$ . To start this problem we will first evaluate

$$\frac{\partial \tilde{f}}{\partial z_i}.$$

Using the chain rule this can be written (and evaluated using the Kronecker delta) as

$$\begin{aligned} \frac{\partial \tilde{f}}{\partial z_i} &= \sum_{j=1}^n \frac{\partial f}{\partial x_j} \frac{\partial x_j}{\partial z_i} = \sum_{j=1}^n \frac{\partial f}{\partial x_j} \frac{\partial}{\partial z_i} \left( \sum_{k=1}^n S_{jk} z_k + s_j \right) = \sum_{j=1}^n \frac{\partial f}{\partial x_j} \left( \sum_{k=1}^n S_{jk} \delta_{ik} \right) \\ &= \sum_{j=1}^n \frac{\partial f}{\partial x_j} S_{ji} = \sum_{j=1}^n \frac{\partial f}{\partial x_j} (S^T)_{ij} = (S^T \nabla f)_i, \end{aligned} \quad (1)$$

where the notation in the last term of the last line means the  $i$ th component of the vector  $S^T \nabla f$ . As a vector equation we have shown

$$\nabla \tilde{f} = S^T \nabla f.$$

We now seek to evaluate  $\nabla^2 \tilde{f}(z)$ . From Equation 1 above we have

$$\begin{aligned} \frac{\partial^2 \tilde{f}}{\partial z_i \partial z_k} &= \sum_{j=1}^n (S^T)_{ij} \frac{\partial}{\partial z_k} \left( \frac{\partial f}{\partial x_j} \right) \\ &= \sum_{j=1}^n (S^T)_{ij} \sum_{l=1}^n \frac{\partial^2 f}{\partial x_l \partial x_j} \frac{\partial x_l}{\partial z_k} \\ &= \sum_{j=1}^n (S^T)_{ij} \sum_{l=1}^n \frac{\partial^2 f}{\partial x_l \partial x_j} \frac{\partial}{\partial z_k} \left( \sum_{q=1}^n S_{lq} z_q + s_l \right) \\ &= \sum_{j=1}^n (S^T)_{ij} \sum_{l=1}^n \frac{\partial^2 f}{\partial x_l \partial x_j} \left( \sum_{q=1}^n S_{lq} \delta_{qk} \right) \\ &= \sum_{j=1}^n (S^T)_{ij} \sum_{l=1}^n \frac{\partial^2 f}{\partial x_l \partial x_j} S_{lk} \\ &= \sum_{j=1}^n (S^T)_{ij} (\nabla^2 f S)_{jk}. \end{aligned}$$

Here  $(\nabla^2 f S)_{jk}$  is the  $jk$ th element of the matrix product  $\nabla^2 f S$ . Note that the above sum is the  $ik$ th element of the matrix product  $S^T \nabla^2 f S$  and we have the identity we were trying to show.

## Problem 2.10

In terms of looking for the minimum of  $f$  as a function of  $x$  following the prescription for search directions in line search methods we will iterate

$$\begin{aligned} s_k &= x_{k+1} - x_k \\ y_k &= \nabla f_{k+1} - \nabla f_k, \end{aligned}$$

with  $B_{k+1}$  from the books Eq. 2.17 or Eq. 2.18 depending on the method used and starting from an initial value  $x_0$ .

In terms of looking for the minimum of  $\tilde{f}(z)$  if we start from a point  $z_0$  and enforce that  $x_k = Sz_k + s$  for all  $k \geq 0$  then the above two equations in terms of  $z_k$  become

$$\begin{aligned} s_k &= Sz_{k+1} - Sz_k = S(z_{k+1} - z_k) \equiv S\tilde{s}_k \\ y_k &= S^{-T}\nabla\tilde{f}_{k+1} - S^{-T}\nabla\tilde{f}_k = S^{-T}(\nabla\tilde{f}_{k+1} - \nabla\tilde{f}_k) \equiv S^{-T}\tilde{y}_k. \end{aligned}$$

Here I have used the results from the previous problem and defined the vector  $\tilde{s}_k$  and  $\tilde{y}_k$ .

We now ask how does the update equation for  $B_{k+1}$  look in terms of these variables  $\tilde{s}_k$  and  $\tilde{y}_k$ . For Eq. 2.17 we have

$$B_{k+1} = B_k + \frac{(y_k - B_k s_k)(y_k - B_k s_k)^T}{(y_k - B_k s_k)^T s_k}, \quad (2)$$

becoming

$$B_{k+1} = B_k + \frac{(S^{-T}\tilde{y}_k - B_k S\tilde{s}_k)(S^{-T}\tilde{y}_k - B_k S\tilde{s}_k)^T}{(S^{-T}\tilde{y}_k - B_k S\tilde{s}_k)^T S\tilde{s}_k},$$

or

$$B_{k+1} = B_k + \frac{S^{-T}(\tilde{y}_k - S^T B_k S\tilde{s}_k)(\tilde{y}_k - S^T B_k S\tilde{s}_k)^T S^{-1}}{(\tilde{y}_k - S^T B_k S\tilde{s}_k)^T \tilde{s}_k},$$

or finally by pre-multiplying by  $S^T$  and post-multiplying by  $S$  we get

$$S^T B_{k+1} S = S^T B_k S + \frac{(\tilde{y}_k - S^T B_k S\tilde{s}_k)(\tilde{y}_k - S^T B_k S\tilde{s}_k)^T}{(\tilde{y}_k - S^T B_k S\tilde{s}_k)^T \tilde{s}_k}.$$

Note that this is Equation 2 in terms of the variable  $z$  when we replace  $B_k \rightarrow S^T B_k S$ .

We now ask how does the update equation for  $B_{k+1}$  look in terms of these variables  $\tilde{s}_k$  and  $\tilde{y}_k$ . For Eq. 2.18 we have

$$B_{k+1} = B_k - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \frac{y_k y_k^T}{y_k^T s_k}, \quad (3)$$

which becomes in this case

$$B_{k+1} = B_k - \frac{B_k S\tilde{s}_k \tilde{s}_k^T S^T B_k}{\tilde{s}_k^T S^T B_k S\tilde{s}_k} + \frac{S^{-T}\tilde{y}_k \tilde{y}_k^T S^{-1}}{\tilde{y}_k^T S^{-1} S\tilde{s}_k},$$

or

$$B_{k+1} = B_k - S^{-T} \left[ \frac{(S^T B_k S)\tilde{s}_k \tilde{s}_k^T (S^T B_k S)}{\tilde{s}_k^T (S^T B_k S)\tilde{s}_k} - \frac{\tilde{y}_k \tilde{y}_k^T}{\tilde{y}_k^T \tilde{s}_k} \right] S^{-1}.$$

Again if we pre-multiplying by  $S^T$  and post-multiplying by  $S$  we get

$$S^T B_{k+1} S = S^T B_k S - \frac{(S^T B_k S) \tilde{s}_k \tilde{s}_k^T (S^T B_k S)}{\tilde{s}_k^T (S^T B_k S) \tilde{s}_k} + \frac{\tilde{y}_k \tilde{y}_k^T}{\tilde{y}_k^T \tilde{s}_k}.$$

Note that this is Equation 3 in terms of the variable  $z$  when we replace  $B_k \rightarrow S^T B_k S$ .

### Problem 2.11

I was not fully sure I understood what this problem was asking but it is easy to imagine a situation where  $f(x)$  is poorly scaled and  $\nabla^2 f$  is ill-conditioned. For example for a scaling of  $\Delta x$  in the  $x$  direction and of  $\Delta y$  in the  $y$  direction we will normally have

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2} &\approx \frac{1}{\Delta x^2} \\ \frac{\partial^2 f}{\partial y^2} &\approx \frac{1}{\Delta y^2} \\ \frac{\partial^2 f}{\partial x \partial y} &\approx \frac{1}{\Delta x \Delta y}. \end{aligned}$$

If we construct  $f$  such that  $\Delta y \gg \Delta x$  and to make our life simpler take  $\frac{\partial^2 f}{\partial x \partial y} = 0$  then we will have

$$\nabla^2 f(x^*) = \begin{bmatrix} \frac{1}{\Delta x^2} & 0 \\ 0 & \frac{1}{\Delta y^2} \end{bmatrix}.$$

If we have  $\Delta y = 10^p \Delta x$  (due to the poor scaling) this matrix is

$$\nabla^2 f(x^*) = \frac{1}{\Delta x^2} \begin{bmatrix} 1 & 0 \\ 0 & 10^{-2p} \end{bmatrix}.$$

As the condition number of a matrix is the ratio of the maximum eigenvalue to the minimum eigenvalue from the above we see that

$$\kappa(\nabla^2 f(x^*)) = \frac{1}{10^{-2p}} = 10^{2p},$$

which can be quite large. A specific example of a function  $f$  that has the above properties is  $f(x_1, x_2) = 10^9 x_1^2 + x_2^2$  which following the arguments above can be shown to have  $\kappa(\nabla^2 f(x^*)) = 10^9$ .

### Problem 2.12

To be Q-linearly convergent we must have

$$\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} \leq r, \quad (4)$$

for  $0 < r < 1$  and all  $k$  sufficiently large. For this sequence the limit would be  $x^* = 0$  and

$$\frac{\|x_{k+1}\|}{\|x_k\|} = \frac{\frac{1}{k+1}}{\frac{1}{k}} = \frac{k}{k+1} = \frac{1}{1 + \frac{1}{k}} > 1,$$

for all  $k$ . Thus  $x_k$  cannot be Q-linearly convergent.

**Problem 2.13**

For this sequence the limit would be  $x^* = 1$  and

$$\frac{\|x_{k+1} - 1\|}{\|x_k - 1\|^2} = \frac{0.5^{2^{k+1}}}{(0.5^{2^k})^2} = \frac{0.5^{2^{k+1}}}{0.5^{2^{k+1}}} = 1 \leq M,$$

for any  $M \geq 1$  and for all  $k$ . Thus  $x_k$  is Q-quadratically convergent to one.

**Problem 2.14**

For this sequence the limit would be  $x^* = 0$  and we need to consider

$$\frac{\|x_{k+1}\|}{\|x_k\|} = \frac{\frac{1}{(k+1)!}}{\frac{1}{k!}} = \frac{1}{k+1} \leq \frac{1}{2},$$

for all  $k > 1$ . This means that  $x_k$  is Q-linearly convergent to zero.

Next consider

$$\frac{\|x_{k+1}\|}{\|x_k\|^2} = \frac{\frac{1}{(k+1)!}}{\left(\frac{1}{k!}\right)^2} = \frac{(k!)^2}{(k+1)!} = \frac{k!}{k+1} \rightarrow \infty,$$

as  $k \rightarrow \infty$ . Thus this sequence does *not* converge Q-quadratically to zero.

**Problem 2.15**

For this sequence, the limit would be  $x^* = 0$  and we need to consider

$$\frac{\|x_{k+1}\|}{\|x_k\|}. \tag{5}$$

Now for  $k$  even  $k+1$  is odd and Equation 5 becomes

$$\frac{\|x_k/k\|}{\|x_k\|} = \frac{1}{k} < \frac{1}{2}, \tag{6}$$

for  $k \geq 2$ . For  $k$  odd,  $k+1$  is even so Equation 5 becomes

$$\begin{aligned} \frac{\|x_{k+1}\|}{\|x_k\|} &= \frac{1}{\|x_k\|} \left(\frac{1}{4}\right)^{2^{k+1}} = \frac{k}{\|x_{k-1}\|} \left(\frac{1}{4}\right)^{2^{k+1}} \\ &= \frac{k \left(\frac{1}{4}\right)^{2^{k+1}}}{\left(\frac{1}{4}\right)^{2^{k-1}}} = k \left(\frac{1}{4}\right)^{2^{k+1} - 2^{k-1}} \\ &= k \left(\frac{1}{4}\right)^{2^k(2 - \frac{1}{2})} = k \left(\frac{1}{4}\right)^{\left(\frac{3}{2}\right)2^k} \rightarrow 0, \end{aligned} \tag{7}$$

as  $k \rightarrow \infty$ . We can prove this last fact (the limit) using L' Hospital's rule. Now if the limit of the above sequence is zero it must eventually be smaller than any value specified if we take  $k$  large enough. If we specify the value  $\frac{1}{2}$  then we have shown that

$$\frac{\|x_{k+1}\|}{\|x_k\|} < \frac{1}{2},$$

for  $k$  sufficiently large. This is the condition needed to show Q-linear convergence.

Note that when we combine Equations 6 and 7 we get that

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1}\|}{\|x_k\|} = 0,$$

which is the condition needed for Q-superlinear convergence.

Now to determine if  $x_k$  is Q-quadratically convergent we need to consider the ratio

$$\frac{\|x_{k+1}\|}{\|x_k\|^2}.$$

If  $k$  is even  $k + 1$  is odd and the ratio is given by

$$\frac{1}{k\|x_k\|} = \frac{1}{k\left(\frac{1}{4}\right)^{2k}} \rightarrow \infty,$$

as  $k \rightarrow \infty$  and  $x_k$  cannot be Q-quadratically convergent (since this limit can never be bounded). As another way to see this consider the case where  $k$  is odd then  $k - 1$  and  $k + 1$  are even so we have

$$\frac{\|x_{k+1}\|}{\|x_k\|^2} = \frac{k^2\|x_{k+1}\|}{\|x_{k-1}\|^2} = \frac{k^2\left(\frac{1}{4}\right)^{2k+1}}{\left[\left(\frac{1}{4}\right)^{2k-1}\right]^2} = \frac{k^2\left(\frac{1}{4}\right)^{2k+1}}{\left(\frac{1}{4}\right)^{2k}} = k^2\left(\frac{1}{4}\right) \rightarrow \infty,$$

as  $k \rightarrow \infty$  which again shows that  $x_k$  cannot be Q-quadratically convergent.

Next if  $x_k$  is to be R-quadratically convergent we need to find a positive sequence  $\nu_k$  such that

$$\|x_k - x^*\| \leq \nu_k,$$

for all  $k$  where  $\nu_k$  is Q-quadratically convergent.

To find a sequence  $\nu_k$  note that if  $k$  is odd (then  $k - 1$  is even) and we have

$$\|x_k\| = \frac{\|x_{k-1}\|}{k} < \|x_{k-1}\| = \left(\frac{1}{4}\right)^{2k-1} < \left(\frac{1}{4}\right)^{2k-2},$$

which is true as  $2^{k-2} < 2^{k-1}$  for  $k \geq 2$ . If  $k$  is even then

$$\|x_k\| = \left(\frac{1}{4}\right)^{2k} < \left(\frac{1}{4}\right)^{2k-2}.$$

This motivates taking the definition of  $\nu_k$  as

$$\nu_k = \left(\frac{1}{4}\right)^{2^{k-2}},$$

for  $k \geq 2$ . We now ask if  $\nu_k$  converges Q-quadratically to zero? To answer this we need to study

$$\frac{\|\nu_{k+1}\|}{\|\nu_k\|^2} = \frac{\left(\frac{1}{4}\right)^{2^{k-1}}}{\left[\left(\frac{1}{4}\right)^{2^{k-2}}\right]^2} = \frac{\left(\frac{1}{4}\right)^{2^{k-1}}}{\left(\frac{1}{4}\right)^{2^{k-1}}} = 1.$$

This is certainly bounded. Because of this we have that  $\nu_k$  is Q-quadratically convergent to zero and that  $x_k$  is thus R-quadratically convergent to zero.

# Chapter 5 (Conjugate Gradient Methods)

## Notes on the Text

### Notes on the linear conjugate gradient method

The objective function we seek to minimize is given by

$$\phi(x) = \frac{1}{2}x^T Ax - b^T x. \quad (8)$$

This will be minimized by taking a step from the current best guess at the minimum,  $x_k$ , in the conjugate directions,  $p_k$ , to get a new best guess as

$$x_{k+1} = x_k + \alpha p_k. \quad (9)$$

Note that to find the value of  $\alpha$  to use in the step given in Equation 9, we can select the value of  $\alpha$  that minimizes  $\phi(\alpha) \equiv \phi(x_k + \alpha p_k)$ . To find this value we take the derivative of this expression with respect to  $\alpha$ , set the resulting expression equal to zero, and solve for  $\alpha$ . Since  $\phi(x_k + \alpha p_k)$  can be written as

$$\begin{aligned} \phi(\alpha) &\equiv \frac{1}{2}(x_k + \alpha p_k)^T A(x_k + \alpha p_k) - b^T(x_k + \alpha p_k) \\ &= \frac{1}{2}x_k^T Ax_k + \alpha x_k^T Ap_k + \frac{1}{2}\alpha^2 p_k^T Ap_k - b^T x_k - \alpha b^T p_k \\ &= \phi(x_k) + \frac{1}{2}\alpha^2 p_k^T Ap_k + \alpha(p_k^T Ax_k - p_k^T b)^T \\ &= \phi(x_k) + \frac{1}{2}\alpha^2 p_k^T Ap_k + \alpha r_k^T p_k. \end{aligned}$$

Setting the derivative of this expression equal to zero gives

$$p_k^T Ap_k \alpha + r_k^T p_k = 0.$$

From which when we solve for  $\alpha$  we get the following for the **conjugate direction stepsize**:

$$\alpha = -\frac{r_k^T p_k}{p_k^T Ap_k}, \quad (10)$$

this is equation 5.6 in the book. Notationally we can add a subscript  $k$  to the variable  $\alpha$  as in  $\alpha_k$  to denote that this is the step size taking in moving from  $x_k$  to  $x_{k+1}$ .

Since the residual  $r$  is defined as  $r = Ax - b$  and to get the the  $k + 1$ st minimization estimate  $x_{k+1}$  from  $x_k$  we use Equation 9. If we premultiply this expression by  $A$  and subtract  $b$  from both sides we get

$$Ax_{k+1} - b = Ax_k - b + \alpha_k Ap_k,$$

or in terms of the residuals we get the **residual update equation**:

$$r_{k+1} = r_k + \alpha_k Ap_k, \quad (11)$$

which is the books equation 5.9.

## The initial induction step in the expanding subspace minimization theorem

I found it a bit hard to verify the truth of the *initial* inductive statement used in the proof of the expanding subspace minimization theorem presented in the book. In that proof the initial step requires that one verify  $r_1^T p_0 = 0$ . This can be shown as follows. The first residual  $r_1$  is given by

$$r_1 = \nabla\phi(x_0 + \alpha_0 p_0) = A(x_0 + \alpha_0 p_0) - b.$$

Thus the inner product of  $r_1$  with  $p_0$  gives

$$r_1^T p_0 = (Ax_0 + \alpha_0 A p_0 - b)^T p_0 = x_0^T A p_0 + \alpha_0 p_0^T A p_0 - b^T p_0.$$

If we put in the value of  $\alpha_0 = -\frac{r_0^T p_0}{p_0^T A p_0}$  the above expression becomes

$$\begin{aligned} &= x_0^T A p_0 - \frac{r_0^T p_0}{p_0^T A p_0} p_0^T A p_0 - b^T p_0 \\ &= (x_0^T A - b^T) p_0 - r_0^T p_0 \\ &= r_0^T p - r_0^T p_0 = 0, \end{aligned}$$

as we were to show.

## Notes on the basic properties of the conjugate gradient method

We want to pick a value for  $\beta_k$  such that the *new* expression for  $p_k$  given by

$$p_k = -r_k + \beta_k p_{k-1},$$

to be  $A$  conjugate to the old  $p_{k-1}$ . To enforce this conjugacy, multiply this expression by  $p_{k-1}^T A$  on the left of the expression above where we get

$$p_{k-1}^T A p_k = -p_{k-1}^T A r_k + \beta_k p_{k-1}^T A p_{k-1}.$$

If we take the left-hand-side of this expression equal to zero and solve for  $\beta_k$  we get that

$$\beta_k = \frac{p_{k-1}^T A r_k}{p_{k-1}^T A p_{k-1}}. \quad (12)$$

In this case the new value of  $p_k$  will be  $A$  conjugate to the old value  $p_{k-1}$ .

## Notes on the preliminary version of the conjugate gradient method

The stepsize in the conjugate direction is given by Equation 10. If we use the conjugate update equation

$$p_k = -r_k + \beta_{k-1} p_{k-1}, \quad (13)$$

and the residual prior-conjugate orthogonality condition given by

$$r_k^T p_j = 0 \quad \text{for } 0 \leq j < k, \quad (14)$$

in Equation 13 we get when we multiply by  $r_k^T$  on the left we have

$$r_k^T p_k = -r_k^T r_k + \beta_{k-1} r_k^T p_{k-1} = -r_k^T r_k,$$

Thus using this fact in Equation 10 we have an alternative expression for  $\alpha_k$  given by

$$\alpha_k = \frac{r_k^T r_k}{p_k^T A p_k}. \quad (15)$$

In the preliminary conjugate gradient algorithm the stepsize,  $\beta_{k+1}$ , in the conjugate direction is given by Equation 12. Using the residual update equation  $r_{k+1} = r_k + \alpha_k A p_k$ , to replace  $A p_k$  in the expression for  $\beta_{k+1}$  to get

$$\beta_{k+1} = \frac{\frac{1}{\alpha_k} (r_{k+1}^T (r_{k+1} - r_k))}{\frac{1}{\alpha_k} p_k^T (r_{k+1} - r_k)} = \frac{r_{k+1}^T r_{k+1}}{p_k^T (r_{k+1} - r_k)},$$

since  $r_{k+1}^T r_k = 0$ , by residual-residual orthogonality

$$r_k^T r_i = 0 \quad \text{for } i = 0, 1, \dots, k-1. \quad (16)$$

Using the conjugate update equation  $p_k = -r_k + \beta_{k+1} p_{k-1}$  in the denominator above, we get a new denominator given by

$$(-r_k + \beta_{k+1} p_{k-1})^T (r_{k+1} - r_k).$$

Next using residual prior-conjugate orthogonality Equation 14 or

$$r_k^T p_i = 0 \quad \text{for } i = 0, 1, \dots, k-1,$$

we have

$$\beta_{k+1} = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}, \quad (17)$$

as we wanted to show.

## Notes on the rate of convergence of the conjugate gradient method

To study the convergence of the conjugate gradient method we first argue that the minimization problem we originally posed: that of minimizing  $\phi(x) = \frac{1}{2} x^T A x - b^T x$  over  $x$  is equivalent to the problem of minimizing a norm squared expression, namely  $\|x - x^*\|_A^2$ . If we take the minimum of  $\phi(x)$  to be denoted as  $x^*$  such that  $x^*$  solves  $Ax = b$  we then have the minimum of  $\phi(x)$  at this point is given by

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

To show that these two minimization problems are equivalent consider the expression  $\frac{1}{2}\|x - x^*\|_A^2$ . We have

$$\begin{aligned}\frac{1}{2}\|x - x^*\|_A^2 &= \frac{1}{2}(x - x^*)^T A(x - x^*) \\ &= \frac{1}{2}x^T Ax - x^T Ax^* + \frac{1}{2}x^{*T} Ax^*.\end{aligned}$$

Since  $Ax^* = b$  we have that the above becomes

$$\begin{aligned}\frac{1}{2}\|x - x^*\|_A^2 &= \frac{1}{2}x^T Ax - x^T b + \frac{1}{2}x^{*T} Ax^* \\ &= \phi(x) + \frac{1}{2}x^{*T} Ax^* \\ &= \phi(x) - \frac{1}{2}x^{*T} Ax^* + x^{*T} Ax^* \\ &= \phi(x) - \frac{1}{2}x^{*T} Ax^* + x^{*T} b \\ &= \phi(x) - \phi(x^*),\end{aligned}$$

verifying the books equation 5.27.

To study convergence of the conjugate gradient method we will decompose the difference between our initial guess at the solution denoted as  $x_0$  and the true solution denoted by  $x^*$  or  $x_0 - x^*$  in terms of the eigenvectors  $v_i$  of  $A$  as  $x_0 - x^* = \sum_{i=1}^n \xi_i v_i$ . When we do this we have that the difference between the  $k + 1$ th iteration and  $x^*$  is given by

$$x_{k+1} - x^* = \sum_{i=1}^n (1 + \lambda_i P_k^*(\lambda_i)) \xi_i v_i.$$

Then in terms of the  $A$  norm this distance is given by

$$\begin{aligned}\|x_{k+1} - x^*\|_A &= \left( \sum_{i=1}^n (1 + \lambda_i P_k^*(\lambda_i)) \xi_i v_i \right)^T A \left( \sum_{i=1}^n (1 + \lambda_i P_k^*(\lambda_i)) \xi_i v_i \right) \\ &= \left( \sum_{i=1}^n (1 + \lambda_i P_k^*(\lambda_i)) \xi_i v_i^T \right) A \left( \sum_{i=1}^n (1 + \lambda_i P_k^*(\lambda_i)) \xi_i v_i \right) \\ &= \sum_{i=1}^n \sum_{j=1}^n \xi_i (1 + \lambda_i P_k^*(\lambda_i)) \xi_j (1 + \lambda_j P_k^*(\lambda_j)) v_i^T A v_j.\end{aligned}$$

As  $v_j$  are orthonormal eigenvectors of  $A$  we have  $A v_j = \lambda_j v_j$  and  $v_i^T v_j = 0$  so that the above becomes

$$\sum_{i=1}^n \xi_i^2 (1 + \lambda_i P_k^*(\lambda_i))^2 \lambda_i,$$

as claimed by the book's equation 5.31.

## Notes on the Polak-Ribiere Method

In this subsection of these notes we derive an alternative expression for  $\beta_{k+1}$  that is used in the conjugate-direction update Equation 13 and that is valid for nonlinear optimization problems. Since our state update equation is given by  $x_{k+1} = x_k + \alpha_k p_k$ , the gradient of  $f(x)$  at the new point  $x_{k+1}$  can be computed to second order using Taylor's theorem as

$$\nabla f_{k+1} = \nabla f_k + \alpha_k \bar{G}_k p_k,$$

where  $\bar{G}_k$  is the *average Hessian* over the line segment  $[x_k, x_{k+1}]$ . To be sure that the new conjugate search direction  $p_{k+1}$  derived from the standard conjugate direction update equation:

$$p_{k+1} = -\nabla f_{k+1} + \beta_{k+1} p_k,$$

is conjugate with respect to the average Hessian  $\bar{G}_k$  means that

$$p_{k+1}^T \bar{G}_k p_k = 0,$$

or using the expression for  $p_{k+1}$  this becomes

$$-\nabla f_{k+1}^T \bar{G}_k p_k + \beta_{k+1} p_k^T \bar{G}_k p_k = 0.$$

So this later expression requires that

$$\beta_{k+1} = \frac{\nabla f_{k+1}^T \bar{G}_k p_k}{p_k^T \bar{G}_k p_k}.$$

Recognizing that  $\bar{G}_k p_k = \frac{\nabla f_{k+1} - \nabla f_k}{\alpha_k}$  so  $\beta_{k+1}$  above becomes

$$\beta_{k+1} = \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{(\nabla f_{k+1} - \nabla f_k)^T p_k}, \quad (18)$$

or the Hestenes-Stiefel formula and the books equation 5.45.

## Problem Solutions

### Problem 5.1 (the conjugate gradient algorithm on Hilbert matrices)

See the MATLAB code `prob_1.m` where we call the MATLAB routine `cgsolve.m` to solve for the solution to  $Ax = b$  when  $A$  is the  $n \times n$  Hilbert matrix (generated in MATLAB using the built-in function `hilb.m`) and  $b$  is a vector of all ones and starting with an initial guess at the minimum of  $x_0 = 0$ . We do this for the values of  $n$  suggested in the text and then plot the number of CG iterations needed to drive the “residual” to  $10^{-6}$ . The result of this calculation is presented in Figure 2. The definition of convergence here is taken to be when

$$\frac{\|Ax - b\|}{\|b\|} < 10^{-6}.$$

We see that the number of iterations grows relatively quickly with the dimension of the matrix  $A$ .

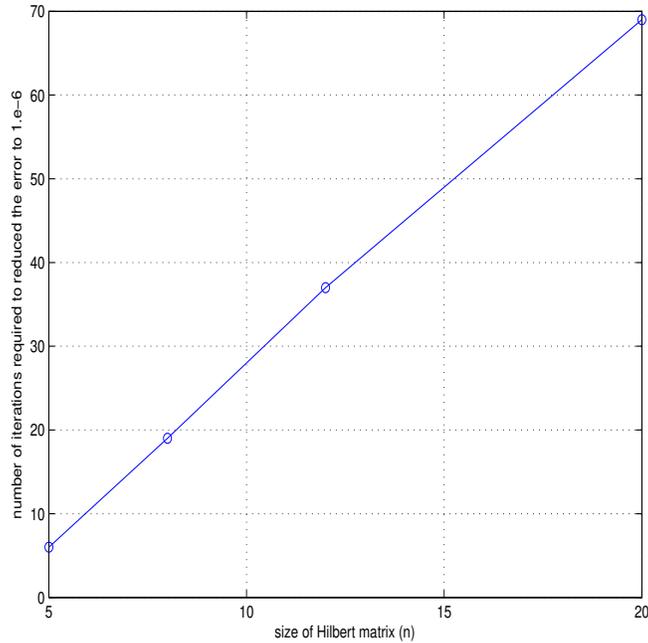


Figure 2: The number of iterations needed for convergence of the CG algorithm when applied to the (classically ill-conditioned) Hilbert matrix.

**Problem 5.2 (if  $p_i$  are conjugate w.r.t.  $A$  then  $p_i$  are linearly independent)**

The book's equation 5.4 is the statement that  $p_i^T A p_j = 0$  for all  $i \neq j$ . To show that the vectors  $p_i$  are linearly independent we begin by assuming that they are not and show that this leads to a contradiction. That  $p_i$  are not linearly independent means that we can find constants  $\alpha_i$  (not all zero) such that

$$\sum \alpha_i p_i = 0.$$

If we premultiply the above summation by the matrix  $A$  we get

$$\sum \alpha_i A p_i = 0.$$

Next premultiply the above by  $p_j^T$  to get

$$\alpha_j p_j^T A p_j = 0,$$

since  $p_j^T A p_i = 0$  for all  $i \neq j$ . Since  $A$  is positive definite the term  $p_j^T A p_j > 0$  we can divide by it and conclude that  $\alpha_j = 0$ . Since the above is true for each value of  $j$  we have that each value of  $\alpha_j$  is zero. This is a contradiction to the assumption that  $p_i$  are not linearly independent thus they must be linearly independent.

**Problem 5.3 (verification of the conjugate direction stepsize  $\alpha_k$ )**

See the discussion around Equation 10 in these notes where the requested expression is derived.

**Problem 5.4 (strongly convex when we step along the conjugate directions  $p_k$ )**

I think this problem is supposed to state that  $f(x)$  is *strongly* convex (as which means that  $f(x)$  has a functional form given by

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

In such a case when we take  $x$  of the form  $x = x_0 + \sigma_0 p_0 + \sigma_1 p_1 + \dots + \sigma_{k-1} p_{k-1}$ , we can write  $f$  as a strongly convex function in the variables  $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_{k-1})^T$ .

**Problem 5.5 (the conjugate directions  $p_i$  span the Krylov subspace)**

We want to show that 5.16 and 5.17 hold for  $k = 1$ . The books equation 5.16 is

$$\text{span}\{r_0, r_1\} = \text{span}\{r_0, Ar_0\}.$$

To show this equivalence we need to show that  $r_1 \in \text{span}\{r_0, Ar_0\}$ . Now  $r_1$  can be written as

$$\begin{aligned} r_1 &= Ax_1 - b \quad \text{with } x_1 \text{ given by} \\ &= A(x_0 + \alpha_0 p_0) - b \quad \text{or} \\ &= Ax_0 - b + \alpha_0 A p_0 \quad \text{or since } r_0 = Ax_0 - b \\ &= r_0 + \alpha_0 A p_0 \quad \text{or since } p_0 = -r_0 \\ &= r_0 - \alpha_0 A r_0, \end{aligned}$$

showing that  $r_1 \in \text{span}\{r_0, Ar_0\}$ , and thus  $\text{span}\{r_0, r_1\} \subset \text{span}\{r_0, Ar_0\}$ . Showing the other direction, that is  $\text{span}\{r_0, Ar_0\} \subset \text{span}\{r_0, r_1\}$  is the same as noting that we can perform the manipulations above in the other direction.

The books equation 5.17 when  $k = 1$  is  $\text{span}\{p_0, p_1\} = \text{span}\{r_0, Ar_0\}$ , since  $p_0 = -r_0$  to show  $\text{span}\{p_0, p_1\} \subset \text{span}\{r_0, Ar_0\}$  we need to show that  $p_1 \in \text{span}\{r_0, Ar_0\}$ . Since

$$\begin{aligned} p_1 &= -r_1 + \beta_1 p_0 \\ &= -(Ax_1 - b) - \beta_1 r_0 \\ &= -(A(x_0 + \alpha_0 p_0) - b) - \beta_1 r_0 \\ &= -r_0 + \alpha_0 A r_0 - \beta_1 r_0. \end{aligned}$$

Again showing the other direction is the same as noting that we can perform the manipulations above in the other direction.

**Problem 5.6 (an alternative form for the conjugate direction stepsize  $\beta_{k+1}$ )**

See the discussion around Equation 17 of these notes where the requested expression is derived.

### Problem 5.7 (the eigensystem of a polynomial expression of a matrix)

Given the eigenvalues  $\lambda_i$  and eigenvectors  $v_i$  of a matrix  $A$  then any polynomial expression of  $A$  say  $P(A)$  has the same eigenvectors with corresponding eigenvalues  $P(\lambda_i)$  as can be seen by simply evaluating  $P(A)v_i$ . This problem is a special case of that result.

### Problem 5.9 (deriving the preconditioned CG algorithm)

For this problem we are to derive the *preconditioned* CG algorithm from the normal CG algorithm. This means that we transform the original problem, that of seeing a solution for the minimum  $\phi(x)$  of

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

by transforming the original  $x$  variable into a “hat” variable  $\hat{x} = Cx$ . In this new space the  $x$  minimization problem above is equivalent to seeking the minimum solution to the following

$$\begin{aligned}\hat{\phi}(\hat{x}) &= \frac{1}{2}(C^{-1}\hat{x})^T A(C^{-1}\hat{x}) - b^T(C^{-1}\hat{x}) \\ &= \frac{1}{2}\hat{x}^T(C^{-T}AC^{-1})\hat{x} - (C^{-T}b)^T.\end{aligned}$$

This later problem we will solve with the CG method where in the standard CG algorithm we take a matrix  $A$  and the vector  $b$  given by

$$\begin{aligned}\hat{A} &= C^{-T}AC^{-1} \\ \hat{b} &= C^{-T}b.\end{aligned}$$

Now given  $\hat{x}_0$  as an initial guess at the minimum of the “hat” problem (note that specifying this is equivalent to specifying a initial guess  $x_0$  for the minimum of  $\phi(x)$ ) and following the standard CG algorithm but using the hated variables, we start to derive the preconditioned CG by setting

$$\begin{aligned}\hat{r}_0 &= C^{-T}AC^{-1}\hat{x}_0 - C^{-T}b \\ \hat{p}_0 &= -\hat{r}_0 \\ k &= 0.\end{aligned}$$

With these initial variables set, we then loop while  $\hat{r}_k \neq 0$  and perform the following steps (following algorithm 5.2)

$$\begin{aligned}\hat{\alpha}_k &= \frac{\hat{r}_k^T \hat{r}_k}{\hat{p}_k^T C^{-T}AC^{-1}\hat{p}_k} \\ \hat{x}_{k+1} &= \hat{x}_k + \hat{\alpha}_k \hat{p}_k \\ \hat{r}_{k+1} &= \hat{r}_k + \hat{\alpha}_k C^{-T}AC^{-1}\hat{p}_k \\ \hat{\beta}_{k+1} &= \frac{\hat{r}_{k+1}^T \hat{r}_{k+1}}{\hat{r}_k^T \hat{r}_k} \\ \hat{p}_{k+1} &= -\hat{r}_{k+1} + \hat{\beta}_{k+1} \hat{p}_k \\ k &= k + 1.\end{aligned}$$

After performing these iterations our output will be  $\hat{x}_\infty$  or the minimum of the quadratic

$$\hat{\phi}(\hat{x}) = \frac{1}{2}\hat{x}^T(C^{-T}AC^{-1})\hat{x} - (C^{-T}b)\hat{x},$$

but we really want to output  $x_\infty = C^{-1}\hat{x}_\infty$ . To derive an expression that works on  $x$  lets write the above algorithm in terms of the unhatted variables  $x$  and  $r$ . Given an initial guess  $x_0$  at the minimum of  $\phi(x)$  then  $\hat{x}_0 = Cx_0$  is the initial guess at the minimum of  $\hat{\phi}(\hat{x})$ . Note that

$$\hat{r}_0 = C^{-T}AC^{-1}\hat{x}_0 - C^{-T}b,$$

or

$$C^T\hat{r}_0 = Ax_0 - b = r_0,$$

is the residual of the original problem. Thus it looks like the residuals transform between hatted an unhatted variables as

$$\hat{r}_k = C^{-T}r_k. \quad (19)$$

Next note that

$$\hat{p}_0 = -\hat{r}_0 = -C^{-T}r_0 = C^{-T}p_0,$$

it looks like the conjugate directions transform between hatted an unhatted variables in the same way, namely

$$\hat{p}_k = C^{-T}p_k. \quad (20)$$

Using these two simplifications our preconditioned conjugate gradient algorithm becomes in terms of the unhatted variables (recall the unknown variable transforms as  $\hat{x}_k = Cx_k$ )

$$\begin{aligned} \hat{\alpha}_k &= \frac{r_k^T C^{-1} C^{-T} r_k}{p_k^T C^{-1} C^{-T} A C^{-1} C^{-T} p_k} \\ &= \frac{r_k^T (C^T C)^{-1} r_k}{p_k^T (C^T C)^{-1} A (C^T C)^{-1} p_k} \end{aligned} \quad (21)$$

$$\begin{aligned} x_{k+1} &= x_k + \hat{\alpha}_k C^{-1} \hat{p}_k \\ &= x_k + \hat{\alpha}_k C^{-1} C^{-T} p_k = x_k + \hat{\alpha}_k (C^T C)^{-1} p_k \end{aligned} \quad (22)$$

$$\begin{aligned} r_{k+1} &= r_k + \hat{\alpha}_k A C^{-1} C^{-T} p_k \\ &= r_k + \hat{\alpha}_k A (C^T C)^{-1} p_k \end{aligned} \quad (23)$$

$$\begin{aligned} \hat{\beta}_{k+1} &= \frac{r_{k+1}^T C^{-1} C^{-T} r_{k+1}}{r_k^T C^{-1} C^{-T} r_k} \\ &= \frac{r_{k+1}^T (C^T C)^{-1} r_{k+1}}{r_k^T (C^T C)^{-1} r_k} \end{aligned} \quad (24)$$

$$\begin{aligned} p_{k+1} &= -r_{k+1} + \hat{\beta}_{k+1} p_k \\ k &= k + 1. \end{aligned} \quad (25)$$

In the above expressions on each line, we first made the hat to unhat substitution and then on the subsequent line simplified the resulting expression. Next to simplify these expressions further we introduce two new variables. The first variable,  $y_k$ , is defined by

$$y_k = (C^T C)^{-1} r_k,$$

or the solution  $y_k$  to the linear system  $My_k = r_k$ , where  $M = C^T C$ . The second variable  $z_k$  is defined similarly as

$$z_k = (C^T C)^{-1} p_k,$$

or the solution to the linear system  $Mz_k = p_k$ . To use these variables, as a first step, in Equation 25 above we multiply by  $(C^T C)^{-1}$  on both sides and use the definitions of  $z_k$  and  $y_k$  to get

$$z_{k+1} = -y_{k+1} + \hat{\beta}_{k+1} z_k.$$

our algorithm then becomes

Given  $x_0$ , our initial guess at the minimum of  $\phi(x)$  form  $r_0 = Ax_0 - b$  and solve  $My_0 = r_0$  for  $y_0$ . Next compute  $z_0$  given by

$$z_0 = M^{-1} p_0 = M^{-1}(-r_0) = -M^{-1}(My_0) = -y_0. \quad (26)$$

Next we set  $k = 0$  and iterate the following while  $r_k \neq 0$

$$\begin{aligned} \hat{\alpha}_k &= \frac{r_k^T y_k}{z_k^T A z_k} \\ x_{k+1} &= x_k + \hat{\alpha}_k z_k \\ r_{k+1} &= r_k + \hat{\alpha}_k A z_k \\ \text{solve } My_{k+1} &= r_{k+1} \text{ for } y_{k+1} \\ \hat{\beta}_{k+1} &= \frac{r_{k+1}^T y_{k+1}}{r_k^T y_k} \\ z_{k+1} &= -y_{k+1} + \hat{\beta}_{k+1} z_k \\ k &= k + 1. \end{aligned}$$

Note that this is the same algorithm presented in the book but the book denotes the variable  $z_k$  by the notation  $p_k$  and I think that there is an error in the book's initialization of this routine in that the book states  $p_0 = -r_0$  while I think this expression should be  $p_0 = -y_0$  (in their notation) see Equation 26.

### Problem 5.10 (deriving the modified residual conjugacy condition)

In the transformed “hat” problem to minimize the quadratic  $\hat{\phi}(\hat{x})$  given by

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

see Problem 5.9 on page 20 above, if we define

$$\begin{aligned} \hat{A} &\equiv C^{-T} A C^{-1} \\ \hat{b} &\equiv C^{-T} b. \end{aligned}$$

So that the transformed hat problem has a residual  $\hat{r}$  given by

$$\begin{aligned} \hat{r} &= \hat{A} \hat{x} - \hat{b} \\ &= C^{-T} A C^{-1} C x - C^{-T} b \\ &= C^{-T} (A x - b) = C^{-T} r. \end{aligned}$$

Thus the orthogonality property of successive residuals i.e. the books equation 5.15 for the hat problem which is given by

$$\hat{r}_k^T \hat{r}_i = 0 \quad \text{for } i = 0, 1, 2, \dots, k-1. \quad (27)$$

becomes in terms of the original variables of the problem

$$r_k^T C^{-1} C^{-T} r_j = r_k^T M^{-1} r_j = 0,$$

since  $M = C^T C$  or the modified residual conjugacy condition and is what we were to show.

**Problem 5.11 (the expressions for  $\beta^{\text{PR}}$  and  $\beta^{\text{HS}}$  reduce to  $\beta^{\text{FR}}$ )**

Recall that three expressions suggested for  $\beta$  (the conjugate direction stepsize) are

$$\beta_{k+1}^{\text{PR}} = \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{\nabla f_k^T \nabla f_k}, \quad (28)$$

for the Polak-Ribiere formula,

$$\beta_{k+1}^{\text{HS}} = \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{(\nabla f_{k+1} - \nabla f_k)^T p_k}, \quad (29)$$

for the Hestenes-Stiefel and

$$\beta_{k+1}^{\text{FR}} = \frac{\nabla f_{k+1}^T \nabla f_{k+1}}{\nabla f_k^T \nabla f_k}, \quad (30)$$

for the Fletcher-Reeves expression.

To show that the Polak-Ribere CG stepsize  $\beta^{\text{PR}}$  reduces to the Fletcher Reeves CG stepsize  $\beta^{\text{FR}}$ , under the conditions given in this problem, from the above formulas it is sufficient to show that

$$\nabla f_{k+1}^T \nabla f_k = 0,$$

since they agree on the other terms. Now when  $f(x)$  is a quadratic function  $f(x) = \frac{1}{2}x^T A x - b^T x + c$  for some matrix  $A$ , vector  $b$ , and scalar  $c$  and  $x_{k+1}$  is chosen to be the exact line search minimum then  $\nabla f_{k+1} = r_{k+1}$ . Thus  $\nabla f_{k+1}^T \nabla f_k = r_{k+1}^T r_k = 0$ , by residual-residual orthogonality Equation 16 (the books equation 5.15).

Now to show that Hestenes-Stiefel CG stepsize  $\beta^{\text{HS}}$  reduces to the Fletcher Reeves CG stepsize  $\beta^{\text{FR}}$ , under the conditions given in this problem, from the above formulas it is sufficient to show that  $(\nabla f_{k+1} - \nabla f_k)^T p_k = \nabla f_k^T \nabla f_k$ . Using the results above we have

$$\begin{aligned} (\nabla f_{k+1} - \nabla f_k)^T p_k &= (r_{k+1} - r_k)^T (-r_k + \beta_k p_{k-1}) \\ &= -r_{k+1}^T r_k + \beta_k r_{k+1}^T p_{k-1} + r_k^T r_k - \beta_k r_k^T p_{k-1} \\ &= r_k^T r_k = \nabla f_k^T \nabla f_k. \end{aligned}$$

Where in the above we have used residual prior-conjugate orthogonality given by Equation 14 to show  $r_k^T p_{k-1} = 0$ .