

# Exercises in Optimization (ACM 40990 / ACM41030)

Dr Lennon Ó Náraigh

## Exercises #2

### Exercises #2 - More on Line-search Methods

1. In the notes (Chapter 6), it is shown that the Newton method satisfies:

$$\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 \leq C\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 \text{ whenever } \|\mathbf{x}_k - \mathbf{x}_*\|_2 < \delta.$$

If we choose

$$\|\mathbf{x}_0 - \mathbf{x}_*\|_2 < \delta, \text{ and } \|\mathbf{x}_0 - \mathbf{x}_*\|_2 < \frac{1}{2C},$$

then

$$\frac{\|\mathbf{x}_1 - \mathbf{x}_*\|_2}{\|\mathbf{x}_0 - \mathbf{x}_*\|_2} \leq C\|\mathbf{x}_0 - \mathbf{x}_*\|_2 \leq \frac{1}{2}.$$

The aim of this exercise is to show that these inequalities give rise to the following important result:

$$\frac{\|\mathbf{x}_k - \mathbf{x}_*\|_2}{\|\mathbf{x}_0 - \mathbf{x}_*\|_2} \leq \frac{1}{2^{2^k - 1}}. \quad (1)$$

The proof of the inequality (1) can be obtained by the following sequence of steps:

- (a) Write down inequalities for

$$\begin{aligned} \|\mathbf{x}_1 - \mathbf{x}_*\|_2 &\leq \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2, & \|\mathbf{x}_2 - \mathbf{x}_*\|_2 &\leq \frac{1}{2} \|\mathbf{x}_1 - \mathbf{x}_*\|_2, \\ & & \|\mathbf{x}_3 - \mathbf{x}_*\|_2 &\leq \frac{1}{2} \|\mathbf{x}_2 - \mathbf{x}_*\|_2. \end{aligned}$$

- (b) Hence, guess that the general term satisfies

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \frac{1}{2^{p_k}} \|\mathbf{x}_0 - \mathbf{x}_*\|_2,$$

where

$$p_k = 2p_{k-1} + 1. \quad (2)$$

- (c) Equation (2) is a first-order **difference equation** with general solution  $p_k = A + B\lambda^n$ , where  $A$ ,  $B$ , and  $\lambda$  are constants to be determined. Hence, show that  $p_k$  satisfies:

$$p_k = 2^k - 1, \quad k > 1$$

with  $p_1 = 1$ .

- (d) Conclude that

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \frac{1}{2^{2^k-1}} \|\mathbf{x}_0 - \mathbf{x}_*\|_2,$$

and hence,

$$\lim_{k \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}_*\|_2 = 0.$$

For part (a) we have:

$$\begin{aligned} \|\mathbf{x}_1 - \mathbf{x}_*\|_2 &\leq C \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2, \\ &\leq (C \|\mathbf{x}_0 - \mathbf{x}_*\|_2) \|\mathbf{x}_0 - \mathbf{x}_*\|_2, \\ &\leq \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2. \end{aligned}$$

Also,

$$\begin{aligned} \|\mathbf{x}_2 - \mathbf{x}_*\|_2 &\leq C \|\mathbf{x}_1 - \mathbf{x}_*\|_2^2, \\ &\leq C \left( \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2 \right) \left( \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2 \right), \\ &\leq \frac{1}{2^2} (C \|\mathbf{x}_0 - \mathbf{x}_*\|_2) \|\mathbf{x}_0 - \mathbf{x}_*\|_2, \\ &\leq \frac{1}{2^3} \|\mathbf{x}_0 - \mathbf{x}_*\|_2. \end{aligned}$$

Finally,

$$\begin{aligned} \|\mathbf{x}_3 - \mathbf{x}_*\|_2 &\leq C \|\mathbf{x}_2 - \mathbf{x}_*\|_2^2, \\ &\leq C \left( \frac{1}{2^3} \|\mathbf{x}_0 - \mathbf{x}_*\|_2 \right) \left( \frac{1}{2^3} \|\mathbf{x}_0 - \mathbf{x}_*\|_2 \right), \\ &\leq \frac{1}{2^6} (C \|\mathbf{x}_0 - \mathbf{x}_*\|_2) \|\mathbf{x}_0 - \mathbf{x}_*\|_2, \\ &\leq \frac{1}{2^7} \|\mathbf{x}_0 - \mathbf{x}_*\|_2. \end{aligned}$$

For part (b), we guess the pattern,  $p_k$  is the power of two in the general term, and we have:

$$\begin{aligned} p_1 &= 1, \\ p_2 &= 3, \\ p_3 &= 7, \end{aligned}$$

so we guess:

$$p_k = 2p_{k-1} + 1 \tag{3}$$

for  $k > 1$ , and  $p_k = 1$  for  $k = 1$ .

For part (c), we substitute the trial solution  $p_k = A + B\lambda^k$  into the difference equation (3). We have:

$$A + B\lambda^k = 2(A + B\lambda^{k-1}) + 1.$$

We equate coefficients to get  $A = -1$  and  $\lambda = 2$ , hence

$$p_k = B(2^k) - 1.$$

We have  $p_1 = 1$ , hence  $B = 1$  also, hence:

$$p_k = 2^k - 1,$$

hence (part (d)),

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \frac{1}{2^{2^k - 1}} \|\mathbf{x}_0 - \mathbf{x}_*\|_2,$$

hence

$$\lim_{k \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}_*\|_2 = 0,$$

as required.

2. Show that if  $0 < c_2 < c_1 < 1$ , there may be no step lengths that satisfy the Strong Wolfe conditions.

Hint: Consider the quadratic function

$$\phi(\alpha) = a + b\alpha + c\alpha^2,$$

where  $b < 0$  and  $c > 0$ .

We compute  $\phi(0) = a$  and  $\phi'(0) = b$ . SW1 requires  $\phi(\alpha) \leq a + c_1 b\alpha$ , hence:

$$b\alpha + c\alpha^2 \leq c_1 b\alpha,$$

hence  $\alpha = 0$  or

$$\alpha = -b(1 - c_1)/c.$$

Identify  $\alpha_1 = |b|(1 - c_1)/c$  with  $0 < c_1 < 1$ . Thus, for SWC1 to hold we require:

$$0 \leq \alpha \leq \alpha_1 = \frac{|b|(1 - c_1)}{c}.$$

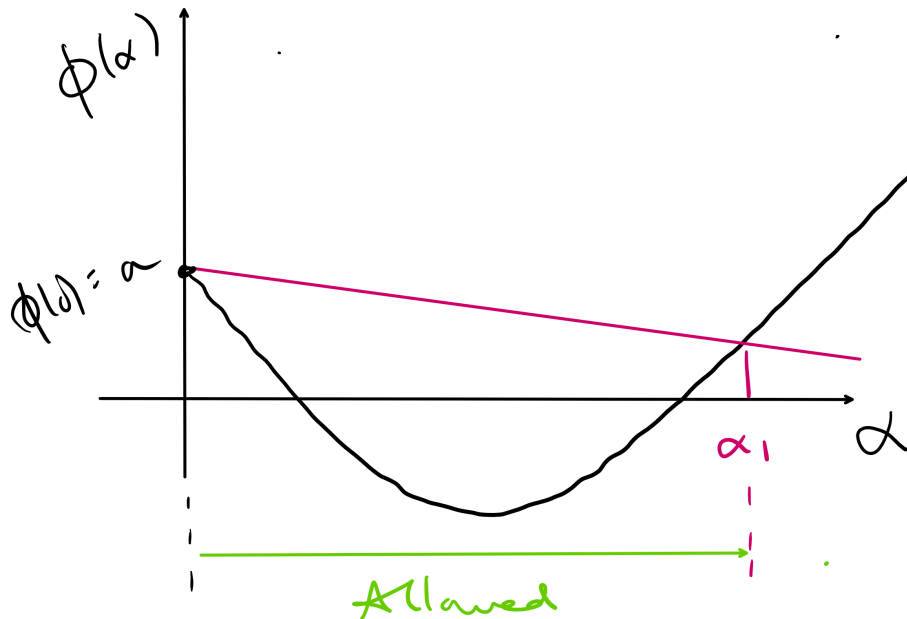


Figure 1: The idea behind SWC1 in the case of a quadratic function  $\phi(\alpha) = a + b\alpha + c\alpha^2$

For SWC2 to hold, we require  $|\phi'(\alpha)| \leq c_2|b|$ , hence  $|b + 2c\alpha| \leq c_2|b|$ . This is a quadratic inequality in disguise:

$$b^2 + 4bc\alpha + 4c^2\alpha^2 \leq c_2^2b^2.$$

We have:

$$4c^2\alpha^2 + 4bc\alpha + b^2(1 - c_2^2) = 0.$$

Hence:

$$\begin{aligned}
 \alpha &= \frac{-4bc \pm \sqrt{16b^2c^2 - 16c^2b^2(1 - c_2^2)}}{8c^2}, \\
 &= \frac{-4bc \pm \sqrt{16c^2b^2c_2^2}}{8c^2}, \\
 &= \frac{-4bc \pm 4|b|cc_2}{8c^2}, \\
 &= \frac{4|b|c \pm 4|b|cc_2}{8c^2}, \\
 &= \frac{|b|(1 \pm c_2)}{2c}, \\
 &= \alpha_{\pm}.
 \end{aligned}$$

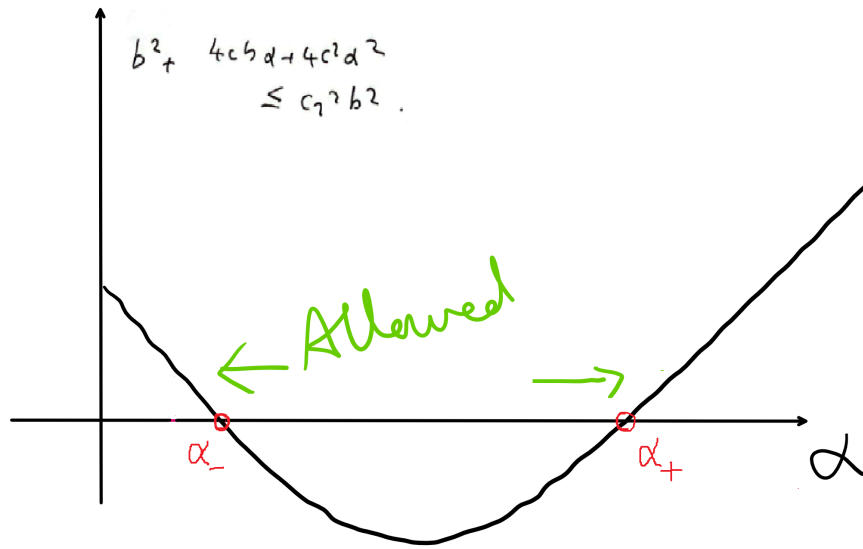


Figure 2: The idea behind SWC2 in the case of a quadratic function  $\phi(\alpha) = a + b\alpha + c\alpha^2$

Referring to Figures 1 and 2, there is a contradiction (and hence, the SWCs don't hold) if

$$\alpha_- > \alpha_1, \quad (4)$$

hence, if:

$$\frac{|b|(1 - c_2)}{2c} > \frac{|b|(1 - c_1)}{c}. \quad (5)$$

Hence, the SWCs do not hold if

$$1 - c_2 > 2 - 2c_1,$$

or

$$2c_1 > c_2 + 1.$$

Refer now to Figure 3. If  $c_1 > c_2$ , then there is the possibility we are in the danger

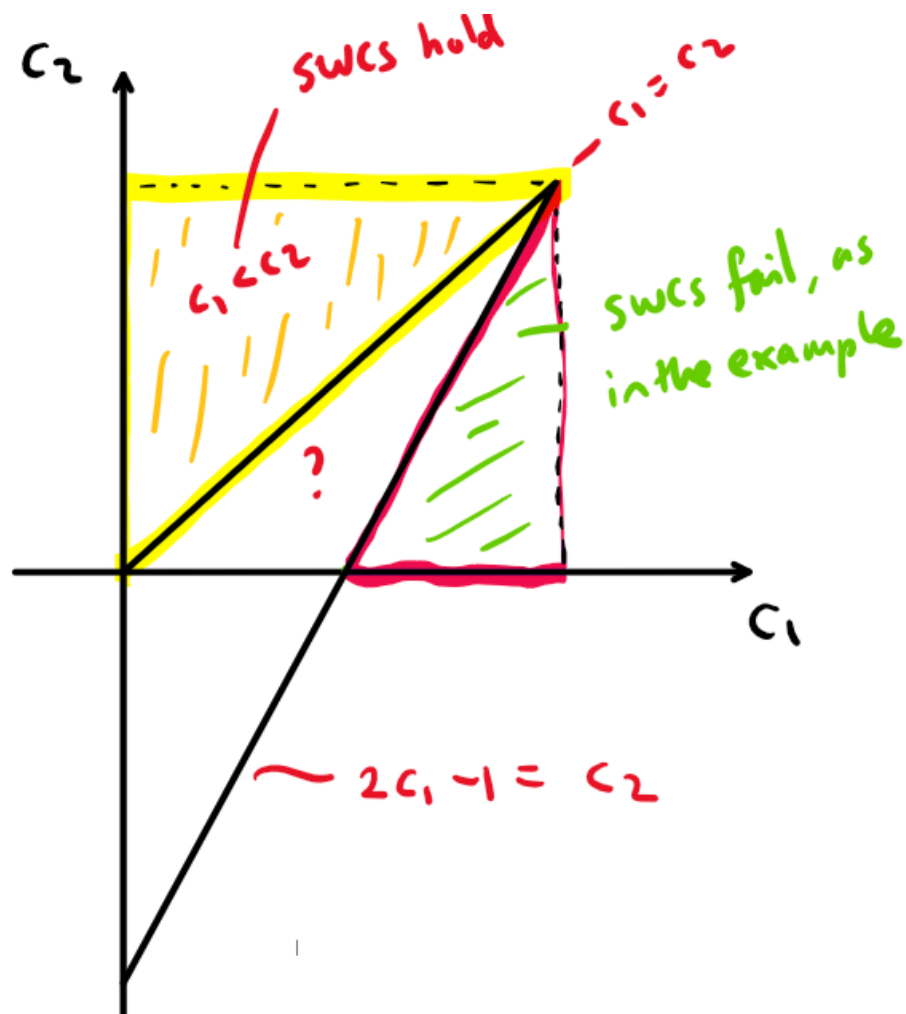


Figure 3: The danger zone (highlighted in red) where the SWCs fail to hold for this example

zone (highlighted in red) in the figure, and hence there is the possibility that the SWCs will fail. Whereas, if  $c_1 < c_2$ , we can't be in the danger zone, and the SWCs will always hold.

3. Consider the one-dimensional function

$$\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{p}_k),$$

where  $\mathbf{p}_k$  is a descent direction – that is,  $\phi'(0) < 0$  – so that our search can be confined to positive values of  $\alpha$ . Find the value that minimizes  $\phi(\alpha)$  in the case where the cost function is quadratic, specifically:

$$f(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{x}, B\mathbf{x} \rangle, \quad (6)$$

where  $\mathbf{a} \in \mathbb{R}^n$  and  $B \in \mathbb{R}^{n \times n}$ .

We have:

$$\begin{aligned} \phi(\alpha) &= f(\mathbf{x}_k + \alpha \mathbf{p}_k), \\ &= \langle \mathbf{a}, \mathbf{x}_k + \alpha \mathbf{p}_k \rangle + \frac{1}{2} \langle \mathbf{x}_k + \alpha \mathbf{p}_k, B\mathbf{x}_k + \alpha B\mathbf{p}_k \rangle, \\ &= \underbrace{\langle \mathbf{a}, \mathbf{x}_k \rangle}_{f(\mathbf{x}_k)} + \alpha \langle \mathbf{a}, \mathbf{p}_k \rangle + \underbrace{\frac{1}{2} \langle \mathbf{x}_k, B\mathbf{x}_k \rangle}_{\frac{1}{2} \langle B\mathbf{x}_k, \mathbf{x}_k \rangle} + \frac{1}{2} \alpha \langle \mathbf{p}_k, B\mathbf{x}_k \rangle + \frac{1}{2} \alpha \langle \mathbf{x}_k, B\mathbf{p}_k \rangle + \frac{1}{2} \alpha^2 \langle \mathbf{p}_k, B\mathbf{p}_k \rangle. \end{aligned}$$

Continuing thus, we have:

$$\begin{aligned} \phi(\alpha) &= \underbrace{f(\mathbf{x}_k)}_{f(\mathbf{x}_k)} + \alpha \langle \mathbf{a}, \mathbf{p}_k \rangle + \frac{1}{2} \alpha \langle \mathbf{p}_k, B\mathbf{x}_k \rangle + \frac{1}{2} \alpha \langle \mathbf{x}_k, B\mathbf{p}_k \rangle + \frac{1}{2} \alpha^2 \langle \mathbf{p}_k, B\mathbf{p}_k \rangle, \\ &= f(\mathbf{x}_k) + \alpha \langle \mathbf{a}, \mathbf{p}_k \rangle + \frac{1}{2} \alpha \langle \mathbf{p}_k, B\mathbf{x}_k \rangle + \frac{1}{2} \alpha \langle B^T \mathbf{x}_k, \mathbf{p}_k \rangle + \frac{1}{2} \alpha^2 \langle \mathbf{p}_k, B\mathbf{p}_k \rangle, \\ &= f(\mathbf{x}_k) + \alpha \langle \mathbf{a}, \mathbf{p}_k \rangle + \frac{1}{2} \alpha \langle \mathbf{p}_k, B\mathbf{x}_k \rangle + \frac{1}{2} \alpha \langle B\mathbf{x}_k, \mathbf{p}_k \rangle + \frac{1}{2} \alpha^2 \langle \mathbf{p}_k, B\mathbf{p}_k \rangle, \\ &= f(\mathbf{x}_k) + \alpha \langle \mathbf{a}, \mathbf{p}_k \rangle + \frac{1}{2} \alpha \langle \mathbf{p}_k, B\mathbf{x}_k \rangle + \frac{1}{2} \alpha \langle B^T \mathbf{x}_k, \mathbf{p}_k \rangle + \frac{1}{2} \alpha^2 \langle \mathbf{p}_k, B\mathbf{p}_k \rangle, \\ &= f(\mathbf{x}_k) + \alpha \langle \mathbf{a}, \mathbf{p}_k \rangle + \alpha \langle \mathbf{p}_k, B\mathbf{x}_k \rangle + \frac{1}{2} \alpha^2 \langle \mathbf{p}_k, B\mathbf{p}_k \rangle, \\ &= f(\mathbf{x}_k) + \alpha \langle \mathbf{p}_k, \mathbf{a} + B\mathbf{x}_k \rangle + \frac{1}{2} \alpha^2 \langle \mathbf{p}_k, B\mathbf{p}_k \rangle. \end{aligned}$$

Hence,

$$\phi'(\alpha) = \langle \mathbf{p}_k, B\mathbf{x}_k + \mathbf{a} \rangle + \alpha \langle \mathbf{p}_k, B\mathbf{p}_k \rangle.$$

When  $\phi'(\alpha) = 0$  we have:

$$\alpha = -\frac{\langle \mathbf{p}_k, B\mathbf{x}_k + \mathbf{a} \rangle}{\langle \mathbf{p}_k, B\mathbf{p}_k \rangle},$$

or

$$\alpha = -\frac{\langle \mathbf{p}_k, \nabla f_k \rangle}{\langle \mathbf{p}_k, B\mathbf{p}_k \rangle},$$

4. Consider the steepest decent method with exact line searches applied to the convex quadratic function in Equation (6).
- (a) Show that if the initial point is such that  $\mathbf{x}_0 - \mathbf{x}_*$  is parallel to an eigenvector of  $B$ , then the steepest descent method will find the solution in one step.
  - (b) Show that the Newton method always converges in exactly one step when the cost function is quadratic, i.e. takes the form (6).

For Part (a) we have:

$$\begin{aligned}\nabla f_0 &= B\mathbf{x}_0 + \mathbf{a}, \\ &= B\mathbf{x}_0 + (-B\mathbf{x}_*), \\ &= B(\mathbf{x}_0 - \mathbf{x}_*).\end{aligned}$$

But  $\mathbf{x}_0 - \mathbf{x}_*$  is parallel to an eigenvector (eigenvalue:  $\lambda$ ), hence:

$$\begin{aligned}\nabla f_0 &= B(\mathbf{x}_0 - \mathbf{x}_*), \\ &= \lambda(\mathbf{x}_0 - \mathbf{x}_*).\end{aligned}$$

From the previous question, we have the value for the exact stepsize, and we are using the SD algorithm, hence  $\mathbf{p}_0 = -\nabla f_0$ . Hence,

$$\begin{aligned}\alpha &= -\frac{\langle \mathbf{p}_0, \nabla f_0 \rangle}{\langle \mathbf{p}_0, B\mathbf{p}_0 \rangle}, \\ &= \frac{\langle \nabla f_0, \nabla f_0 \rangle}{\langle \nabla f_0, B\nabla f_0 \rangle}, \\ &= \frac{\lambda^2 \langle \mathbf{x}_0 - \mathbf{x}_*, \mathbf{x}_0 - \mathbf{x}_* \rangle}{\lambda^2 \langle \mathbf{x}_0 - \mathbf{x}_*, B(\mathbf{x}_0 - \mathbf{x}_*) \rangle}, \\ &= \frac{1}{\lambda}.\end{aligned}$$

Finally, we have the steepest-descent step:

$$\begin{aligned}\mathbf{x}_1 &= \mathbf{x}_0 - \alpha \nabla f_0, \\ &= \mathbf{x}_0 - \frac{1}{\lambda} [\lambda(\mathbf{x}_0 - \mathbf{x}_*)], \\ &= \mathbf{x}_*,\end{aligned}$$

and thus, the SD algorithm converges in one step.

For Part (b) we have:

$$\begin{aligned}\mathbf{x}_1 &= \mathbf{x}_0 - B^{-1} \nabla f_0, \\ &= \mathbf{x}_0 - B^{-1} [B\mathbf{x}_0 + \mathbf{a}], \\ &= -B^{-1} \mathbf{a}, \\ &= \mathbf{x}_*,\end{aligned}$$

and thus, the Newton algorithm converges in one step.

5. Consider the optimization problem,

$$\min f(\mathbf{x}), \quad f(\mathbf{x}) = \langle \mathbf{a}, \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{x}, B\mathbf{x} \rangle,$$

where now  $B$  is a specific  $10 \times 10$  matrix and  $\mathbf{a}$  is a specific  $10 \times 1$  column vector. The numerical values of these arrays can be found in the spreadsheet `OP_10x10.csv`:

- The spreadsheet contains a  $10 \times 1$  array which corresponds to the vector  $\mathbf{a}$ ;
- The spreadsheet contains a  $10 \times 10$  array  $B_0$ .

The array  $B$  is obtained from  $B_0$  by the following sequence of steps:

(i) Symmetrize  $B_0$ :

$$B_0 \rightarrow (B_0 + B_0^T)/2;$$

(ii) Scale  $B_0$ :

$$B_0 \rightarrow B_0 / \max(|B_0|)$$

(iii) Generate a positive-definite matrix:

$$B_0 \rightarrow (B_0^T)B_0.$$

The end result of this sequence of operations is the matrix  $B$ .

Hence,

- Find the minimizer  $\mathbf{x}_*$  numerically, using the steepest-descent and Newton algorithms.
- Why is the convergence so poor in the case of the steepest-descent algorithm?

For part (a), we present sample codes in Matlab, these can be found in an accompanying folder (`OP_10x10`). The convergence is extremely slow in case of the SD method (for the given numerical parameters); results are not presented here. Instead, we go straight over to the Newton method, where we obtain the following results:

```
>> [x,f,1] = add_w()
Cost Function:17.9407
a=1
Cost Function:-225.0464
a=1
Convergence Reached: \|nabla f\|=9.1863e-12

x =

    1.0e+03 *

    0.248481416364386
    1.318605978729388
    0.379081869087112
   -0.390457125647163
   -0.520952833960023
   -0.210071124292290
   -0.201216895062331
   -0.738932520228526
    0.847612993079193
   -0.261164500199232

f =

   -2.250464354514670e+02

1 =

     2
```

Results using  
Newton Method

```
>> temp=load('problem.mat');
B=temp.ans.B;
a=temp.ans.a;

% B=B(1:5,1:5);
% a=a(1:5);

% B is initially a random matrix generated from a file, now I turn it into
% a symmetric positive definite matrix.

B=(B+B.)/2;
B=B/(max(abs(B)));
B=(B.)*B;
>> x_star=-(B\(-1))*a

x_star =

    1.0e+03 *

    0.248481416365052
    1.318605978729388
    0.379081869088129
   -0.390457129482205
   -0.520952833962220
   -0.210071124292290
   -0.201216895062868
   -0.738932520230494
    0.847612993081450
   -0.261164500199958
```

Known Analytical Result

This can be checked, as we know the analytical value of the minimizer in this case,  $\mathbf{x}_* = -B^{-1}\mathbf{a}$ , this is also shown in the figure. The numerical and analytical solutions agree, confirming the correct implementation of the Newton method in this case.

For part (b), we look at the condition number of the matrix  $B$ :

```
>> cond(B)

ans =

    2.746503264292211e+05
```

This is greater than  $10^5$  meaning the maximum and minimum eigenvalues are orders of magnitude apart. From class notes, and for the SD method, we have:

$$\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_B^2 \leq \left(1 - \frac{1}{\kappa(B)}\right) \|\mathbf{x}_k - \mathbf{x}_*\|_B^2,$$

and with  $\kappa(B) > 10^5$ , we have:

$$\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_B^2 \lesssim \|\mathbf{x}_k - \mathbf{x}_*\|_B^2,$$

leading to poor convergence of the SD method in this particular example.